# C++ Parallel Programming Models

## ERAN GILAD

# Nice to meet you!

Day job:

Database internals @



Community activity:

 Conference

 Meetup

 Chit-chat

# Why parallel programming?

1. Performance          2. Hide I/O          3. Responsiveness

Why is it hard?
non-deterministic execution, data races, deadlocks

Having *structure* makes things easier

# Programming model

Model → structure → abstractions

User can reason about correctness
Library can provide useful features

Runtime can optimize scheduling

Crucial in parallel programming!

# C++ parallel programming facilities

condition_variable     schedule          co_return   co_await

thread   thread_local       then        co_yield   execution::par

latch   atomic   lock_guard   future   async   execution::seq

counting_semaphore   mutex   reduce   exclusive_scan

# C++ parallel programming facilities

condition_variable

thread    thread_local

latch    atomic    lock_guard

counting_semaphore    mutex

schedule
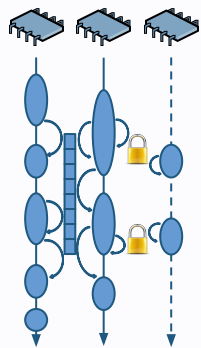
then

future    async

co_return    co_await

co_yield    execution::par

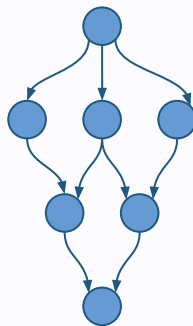execution::seq

reduce    exclusive_scan

Not random facilities, not a single API:
4 parallel programming models

# C++ parallel programming models
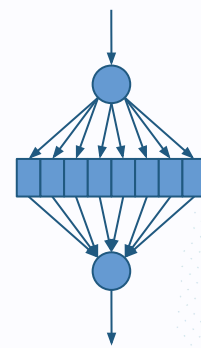


**Unstructured**

thread, mutex, atomics

(C++11)

**Task parallelism**

async, senders

(C++11, C++26)

**Cooperative tasks**

coroutines

(C++20)

**Data parallelism**

parallel algorithms

(C++17)

Not different *abstraction levels* - different program structuring

# Talk outline

Intro

**Unstructured parallelism**

**Task-based parallelism**
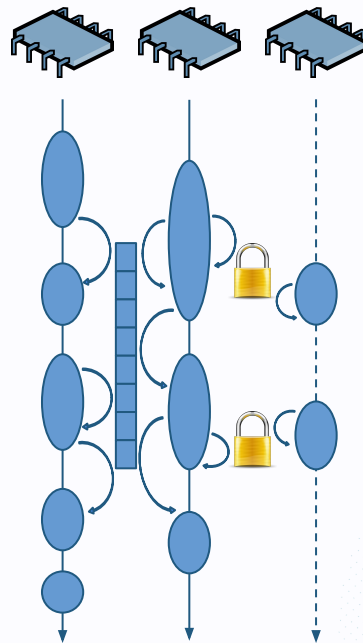
**Cooperative multitasking**

**Data parallelism**

**Mixing models**

*This talk is about the forest,*

*not about the trees*

# Model #1: Unstructured parallelism

- Ad hoc use of parallelization facilities

  - Solve a particular need

- Relies on lower-level abstractions

- Min overhead, max hardware utilization

- Min safety, so requires max proficiency

- Central C++11 feature

# Major components (partial list)

C++11/14

C++20

**Convenience utils**
RAII lock wrappers, call_once, jthread

**Threading and sync primitives**
thread, mutex, condition_variable, semaphore, latch, barrier

**Memory model**
atomic and friends, thread_local

# Exclusive use cases

- Construct higher level facilities
    - E.g., thread pool, spin lock etc.
- Concurrent data structures
    - At least thread safe, usually better if lock free
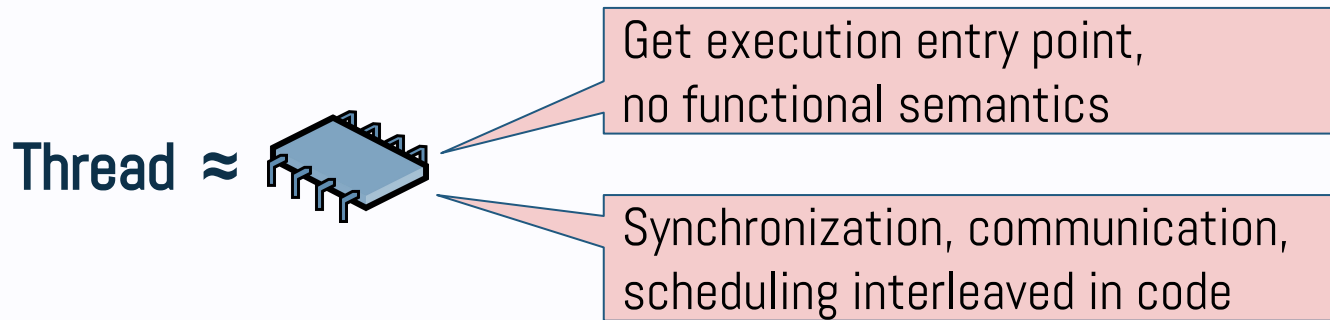- Long running services

# Missing part - safe shared state

- Threads communicate via shared state
    - atomics are too fine-grained
    - locking complete containers is too coarse grained
- We need *concurrent data structures!*
    - Nothing as of C++23
    - RCU and hazard pointers coming up in C++26 🤞
    - Concurrent queues in the works

# Unstructured parallelism pros and cons

- Maximal control

- Usually maximal performance (when done right)

- Complicated memory model, hard to make ideal use

- Data races, deadlocks, non-determinism
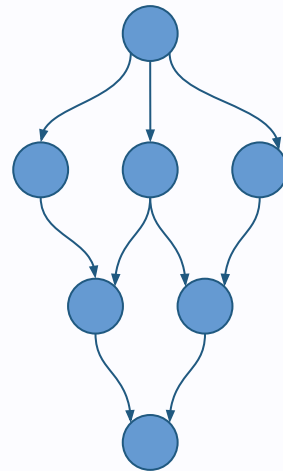
# Thread-level API shortcomings

Thread ≈ [chip illustration]

Get execution entry point,
no functional semantics

Synchronization, communication,
scheduling interleaved in code

Higher abstraction level + clever runtime

⇓

Less work, less bugs, hopefully better performance

# Model #2: Task-based parallelism

- Task: limited computation providing single result
  - Function, lambda, loop iteration etc.
  - *Ideally* has inputs and/or output, no side effects

- Decouples functionality from execution
  - One thread can run many tasks
  - One task can migrate among many threads (?)

# Tasks are asynchronous

Execution:

- Somewhere, sometime
- Creator can proceeds until the outcome is needed

Task results:

- *Future*: A handle to the task outcome
- Or passed directly to a following task

# C++11 tasks - std::async et al.

```cpp
future<int> answer = async(..., []{ return 42; });
do_caller_stuff();
print("found the answer: {}", answer.get());
```

- Spawn tasks using std::async
- Obtain results using std::future
- The runtime assigns tasks to worker threads
  - Yes, the C++ runtime can create and manage threads without user intervention
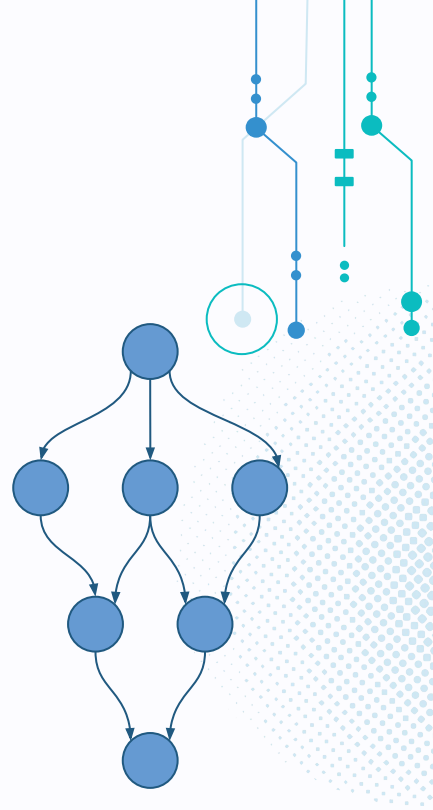
# std::async execution

```
future<int> answer = async(when, []{ return 42; });
do_caller_stuff();
print("found the answer: {}", answer.get());
```

Task execution determined by **when**:

1. **launch::async** - on a new thread ("as if"!)
2. **launch::deferred** - right before the results are needed
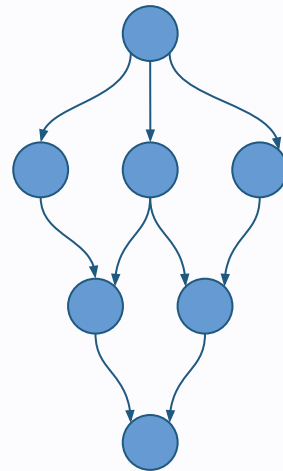3. **async | deferred** - implementation dependent (default)

# std::async tasks limitations

- No composition - follow-up, await many, etc.
- No execution config - e.g., thread pool
- No scheduling options - e.g., yield, stop
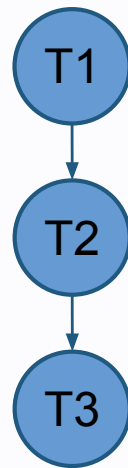- No changes since C++11 🤷

# C++26 Sender/receivers tasks

- Composable by design
  - Chaining, waiting, splitting
- Configurable scheduler
  - Thread pools, coroutines, GPU, etc.
- Support stopping and error reporting
- Coming up in C++26!

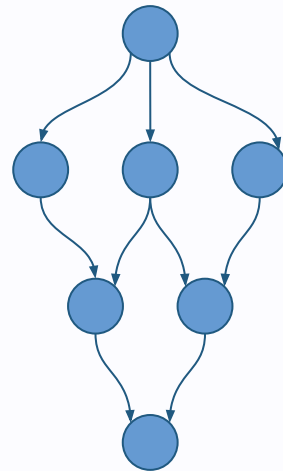# Senders as tasks

```cpp
auto tasks =
  just("world")              // task 1
  | then([](string who) {    // task 2
        return format("Hello {}", who); })
  | then([](string greet) { // task 3
        println("{}!", greet); });

sync_wait(tasks);
```

T1

T2

T3

Senders define a graph of tasks, with input-output channeling

# The future of async C++

- Senders/receivers are extendable
  - Schedulers, thread pools, processor type
- Allow flexible task graphs
  - With clean error handling
- Common infra for async operations
  - Sender performs the async work
  - Receiver called when done

# Model #3: Cooperative Multitasking

- Long running tasks should not block
  - Assuming #tasks > #cores
  - Breaking looping sender tasks can be cumbersome
- Tasks are a user-mode concept
  - The kernel can't help
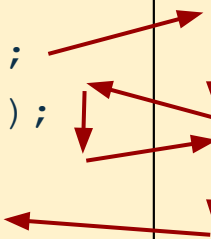
Enter coroutines!

# Coroutines

C++20 coroutine: a function that can be explicitly suspended and resumed by the C++ code

- No operating system involvement
- No assembly required
- Full language (and compiler) support
- Well defined state management and error handling

# Example

```
int main() {                       coro_t coro_foo() {

  coro_t c = coro_foo();             println("Before await");
  println("c suspended");            co_await suspend_always{};
  c.resume();                        println("After await");

}                                  }
```

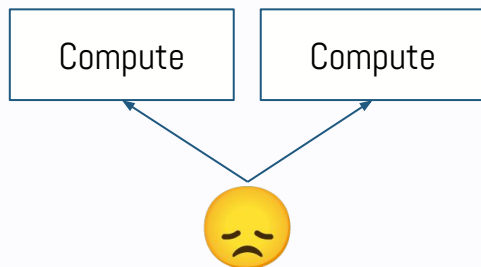coro_t: define coroutine behavior and resume handle

# co_await *expr*

- *expr* should evaluate to an *awaitable object*
- The awaitable is a scheduling point:
  - await_ready() - should the coroutine be suspended
  - await_suspend() - get the suspended coroutine handle
  - await_resume() - after coroutine resumed
- co_await foo():
  - execute foo, evaluate its return value (awaitable)
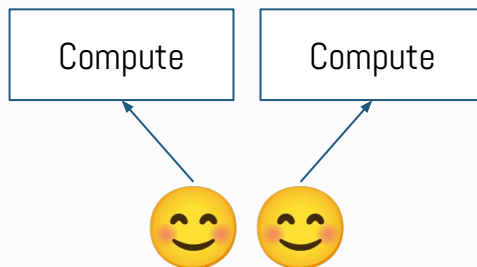  - suspend the current coroutine if needed

# Coroutine scheduling

- co_await returns control to caller or another coroutine
  - It *does not* create any parallelism!
- The program must:
  - Track suspended coroutines
  - Resume ready coroutines
- Requirements:
  - Some source of parallelism (async I/O, worker threads, etc.)
  - Scheduler

# Concurrency and parallelism
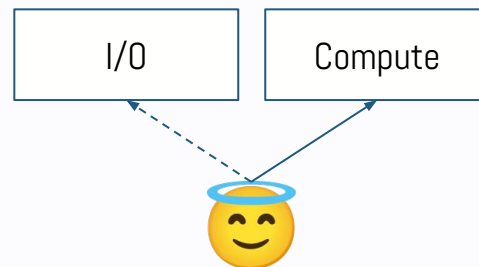
Task concurrency,
no parallelism

| Compute | Compute |
|---------|---------|

Task concurrency,
multi-CPU parallelism

| Compute | Compute |
|---------|---------|

Task concurrency,
CPU/IO parallelism

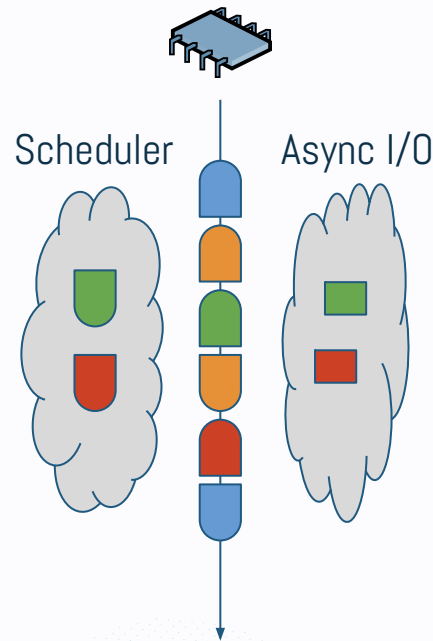| I/O | Compute |
|-----|---------|

# Coroutine-based multitasking example

1. The scheduler starts a coroutine task
2. The task invokes **async** I/O
3. The task **suspends itself**
4. The scheduler switches to another task
5. The async I/O completes
6. The awaiting task marked as ready
7. The scheduler **resumes it** when possible



Scheduler    Async I/O

# Parallel coroutines execution

- Commonly one thread per core
- Suspended tasks can move among cores
  - Can you use thread_local in this scenario?
- Sync primitives must be coroutine-aware
  - co_await coro_lock() instead of std::mutex.lock()
  - A coroutine can't switch if its thread is blocked!
- Expected scheduler complexity - races, load balancing, task affinity, etc.

# Threads vs. tasks

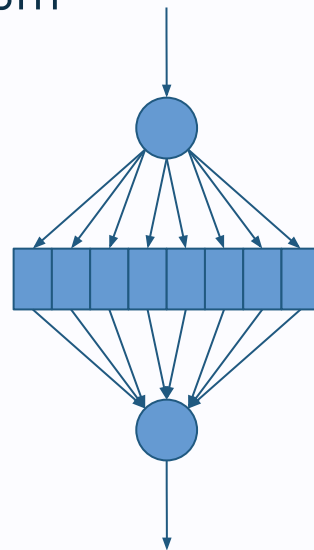| Threads use cases | Tasks use cases |
|---|---|
| Parallelism abstraction (CPU) | Concurrency abstraction (function) |
| Background services/workers | Many short and independent ops |
| Long complex parallel operations | "local"/ad-hoc (async I/O) |
| Equal amount of per-thread work | Dynamic load balancing |

# Fibers / user threads / green threads

- Stackful coroutines
    - A full stack must be allocated for each fiber
- Any function running on a fiber can be suspended
    - Not just the coroutine body
- Decoupled executor (fiber) and work (function)
    - Allows fiber pools etc.
- Requires an "ecosystem" - scheduler and sync objects
- Not supported in C++ (yet?)

# Model #4: Data parallelism

A programming model in which parallelism stems from the individual computations associated with each element in a collection.

- Requires a lot of non-dependent data
- Procedure:
  - Divide element processing among processors
  - Short and simple operation on each element

# C++ Parallel algorithms

```
res = algorithm(
        exec_policy,
        container);
```

Sequential execution context

Declarative API

Accelerate data processing

- High level abstraction
  - Few customization points for user
  - No control over parallelism, scheduling, work distribution etc.
  - Library + runtime can be very efficient
  - When used properly: no data races, deadlocks etc.

# Using parallel algorithms

```cpp
vector<int> v = {1, 2, 3, 5, 11, 20};
int res = reduce(execution::par, v.begin(), v.end());
assert(res == 42);
```

- Most standard algorithms have a parallel overload
  - First parameter: *ExecutionPolicy*
- Complexity requirements more lax
- Implementation isn't specified – can also use GPUs

# sequenced_policy

- Forces execution to take place on the calling thread
- Minor differences from no-policy call
- execution::seq is an *instance of* sequence_policy
  - Algorithms can be overloaded by policy
  - Policy is a *compile time* decision!

```
reduce(execution::seq, ..., ...);
```

# parallel_policy

- Execution on caller or another thread (runtime pool)
- Per thread, semantics are similar to sequenced_policy - unspecified order, no interleaving
- Data races are now possible if multiple operations write to unprotected data

```
reduce(execution::par, ..., ...);
```

# parallel_unsequenced_policy

- A single algorithm thread can now process multiple elements at the same time
    - Operations must not use any locks
- More user restrictions => more library options
    - Vectorization can now be used
    - finer grained scheduling

```
reduce(execution::par_unseq, ..., ...);
```

# unsequenced_policy

- Operations can be interleaved on a single thread
  - Not a multithreaded context
  - But vectorization can still be used!

- C++20 addition

```
reduce(execution::unseq, ..., ...);
```

# Non-standard policies

- Vendor-specific

- Can allow the use of accelerators

    - GPU

    - FPGA

    - ASIC

    - ...

# Parallel algorithms != parallel containers

- C++ separates algorithms from containers

    - Thread-safe containers by default? No - zero-overhead!

- Parallel algorithms can modify data but not structure

    - Unlike sequential algorithms, which can modify both

- Accessing a container processed by a parallel algorithm from another thread? Possible race!

# Mixing models - unstructured context

## Most cores are used:

- E.g., server
- Use coroutines for async I/O
- No resources for parallel algorithms or async compute

## Most cores are unused:

- E.g., background service
- Tasks and parallel algorithms can be used
  - Impact? depends
- Careful with shared state!

Reminder: runtime schedulers aren't aware of user threads

# Mixing models - async/sender tasks context

## Most cores are used:

- Namely, many tasks created
- No point in creating threads or using parallel algorithms
- Awkward context for coroutine execution
  - Even of tasks are coroutines

## Most cores are unused:

- E.g., ad-hoc work, UI worker
- Can use parallel algorithms from tasks
- Creating threads less suitable - spawn tasks instead
- Using sync mechanism (mutex etc.) doesn't fit the model

# Mixing models - coroutine tasks context

## Most cores are used:

- Namely, the scheduler uses many cores
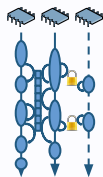- No point in using other models from within coroutines

## Most cores are unused:

- E.g., limited I/O needs
- Parallel algorithms from coroutine will block it
- So will async/sender tasks
- Tasks too short for thread creation

# Mixing models - parallel algorithms

- Within a parallel algorithm (e.g., user lambda):
    - No point in spawning async tasks or creating threads
    - Inappropriate context for coroutines
    - Using sync mechanisms to access external state possible but might kill concurrency

# Summary

**Thank you! Qs?**

**Unstructured**

Low-level building blocks

**Cooperative multitasking**

great for async I/O,
missing scheduler

**Task parallelism**

functional decomposition,
much better with senders

**Data parallelism**

declarative, no executors
control (yet!)

Mixing models hardly works. Parallel programming is hard.