

HSE

Faculty of Computer Science

Конспект углубленно- го курса по C++

AUTHOR

Your Name

COURSE

Advanced C++
Fall 2023

Оглавление

I Лекция 10 – Advanced thread	3
1 Anatomy of a Thread: Cost, Kernel & Scheduling	4
1.1 Физическая структура потока	4
1.2 Стоимость создания потока	5
1.2.1 Бенчмарк: Thread vs Function	5
1.3 Memory Overhead: Стек и Виртуальная память	6
1.4 Планировщик Linux (CFS)	6
1.4.1 Механика работы CFS	6
1.4.2 Квант времени и Переключение контекста	6
1.5 Антипаттерн: Thread per Request	7
1.6 Проблемы стандартных абстракций C++	7
1.6.1 std::execution::par (C++17)	7
1.6.2 std::async и std::future	8
2 Building a Production-Grade ThreadPool	9
2.1 Фундамент: Блокирующая очередь (Blocking Queue)	9
2.1.1 Реализация UnboundedBlockingQueue	9
2.1.2 Разбор механики синхронизации	10
2.2 Архитектура ThreadPool	11
2.2.1 Интерфейс задач	11
2.2.2 Полная реализация ThreadPool	12
2.3 Анализ Corner Cases и Ошибок	13
2.3.1 1. Исключения в воркерах	13
2.3.2 2. Порядок остановки (Destruction Order)	13
2.3.3 3. Невозможность принудительного убийства (Thread Cancellation)	13
2.3.4 4. Data Races при проверке состояния	14
3 User-Space Concurrency: Implementing Fibers	15
3.1 Кооперативная многозадачность	15
3.2 Архитектура Файбера	15
3.2.1 Инфраструктура: Стек	16
3.3 Deep Dive: Assembly & Calling Conventions	16
3.3.1 Магия переключения: SaveContext и JumpContext	17
3.4 Трамплин и Запуск файбера	18
3.5 Планировщик (Scheduler)	19
3.5.1 Intrusive List	19
3.5.2 Цикл планирования	19
3.6 Примитивы синхронизации: Yield и Suspend	20
3.6.1 Yield (Добровольная уступка)	20
3.6.2 Suspend (Блокировка)	20
4 Linux Low-Level Sync: The Futex	22
4.1 Философия Futex	22
4.2 Системный вызов futex()	23
4.2.1 FUTEX_WAIT	23

4.4 Внутри ядра: Kernel Wait Queues	25
5 Lock-Free Programming: Atomic Foundations	26
5.1 Определения и Гарантии	26
5.2 Атомики в C++	26
5.2.1 Compare-And-Swap (CAS)	27
5.3 Паттерн CAS-Loop	27
5.3.1 Strong vs Weak CAS	28
5.4 Пример: Lock-Free Stack (Treiber Stack)	28
5.4.1 Операция Push	28
5.4.2 Операция Pop и утечки памяти	29
5.5 Проблема безопасного удаления памяти	29
5.6 Спинлок (Spinlock) на атомиках	30
6 Advanced Lock-Free: ABA & Memory Reclamation	31
6.1 Проблема АВА	31
6.1.1 Анатомия катастрофы (Сценарий на Стеке)	31
6.2 Решение 1: Tagged Pointers (Указатели с версиями)	32
6.2.1 Реализация на x86_64	32
6.3 Решение 2: Hazard Pointers (Опасные указатели)	33
6.3.1 Алгоритм работы с HP в Pop	33
6.3.2 RetireNode и Scan	34
6.4 Решение 3: RCU (Read-Copy-Update)	34
6.5 Практика: Lock-Free Очередь (Queue)	35
6.5.1 Проблемы двух указателей	35
6.5.2 Dummy Node (Фиктивный узел)	35

Часть I

Лекция 10 – Advanced thread

Глава 1

Anatomy of a Thread: Cost, Kernel & Scheduling

Многопоточность в C++ – это не просто абстракция `std::thread`. За каждым объектом стандартной библиотеки скрывается сложный механизм операционной системы, имеющий свою цену в тактах процессора и байтах оперативной памяти. Понимание физического устройства потока, алгоритмов планировщика ядра (OS Scheduler) и накладных расходов на переключение контекста является фундаментом для построения высоконагруженных систем.

В этой главе мы деконструируем понятие "поток", разберем стоимость его создания и эксплуатации, а также проанализируем архитектурные ошибки, возникающие при наивном использовании примитивов стандарта.

Физическая структура потока

Контекст исполнения (Execution Context)

Поток (thread) физически представляет собой совокупность состояния регистров процессора и выделенного региона памяти под стек. В многопоточной среде все потоки одного процесса разделяют общее адресное пространство (кучу, сегмент кода, глобальные переменные), но имеют изолированный стек и набор регистров.

С точки зрения ядра Linux, поток – это "легковесный процесс" (LWP – Light Weight Process). Для планировщика нет принципиальной разницы между процессом и потоком, кроме того факта, что потоки делят таблицу страниц памяти (Page Table).

Минимальный набор данных, определяющий поток, включает:

1. **Регистры общего назначения (GPR):** RAX, RBX, RCX и т.д. в архитектуре x86_64.
2. **Указатель инструкций (Instruction Pointer, RIP):** Адрес текущей исполняемой команды.
3. **Указатель стека (Stack Pointer, RSP):** Адрес вершины стека данного потока.
4. **Стек (Stack):** Область памяти для локальных переменных, адресов возврата и аргументов функций.

При переключении контекста ядро ОС обязано сохранить текущее состояние регистров уходящего потока в специальную структуру в Kernel Space и загрузить значения регистров приходящего потока.

Стоимость создания потока

Создание потока – операция дорогая. В отличие от вызова функции, который занимает на носекунды (просто `call` и `push` в стек), создание потока (`std::thread`) требует системного вызова (`clone` в Linux), выделения структур ядра и настройки виртуальной памяти.

Эмпирическая оценка времени создания потока на современном железе составляет **от 5 до 50 микросекунд**. Это на несколько порядков медленнее, чем создание объекта в куче.

Бенчмарк: Thread vs Function

Сравним стоимость выполнения пустой работы через вызов функции и через создание нового потока.

```

1 #include <iostream>
2 #include <thread>
3 #include <vector>
4 #include <chrono>
5
6 void dummy_work() {
7     // Имитация минимальной работы
8     volatile int x = 0;
9     x++;
10 }
11
12 int main() {
13     const int N = 10000;
14
15     auto start = std::chrono::high_resolution_clock::now();
16     for (int i = 0; i < N; ++i) {
17         dummy_work();
18     }
19     auto end = std::chrono::high_resolution_clock::now();
20     std::cout << "Function calls: "
21             << std::chrono::duration_cast<std::chrono::microseconds>(end -
22             & start).count()
23             << " us\n";
24
25     start = std::chrono::high_resolution_clock::now();
26     for (int i = 0; i < N; ++i) {
27         std::thread t(dummy_work);
28         t.join();
29     }
30     end = std::chrono::high_resolution_clock::now();
31     std::cout << "Thread creation: "
32             << std::chrono::duration_cast<std::chrono::microseconds>(end -
33             & start).count()
34             << " us\n";
35 }
```

Результаты на типичном сервере показывают, что создание потоков занимает в тысячи раз больше времени. Если ваша задача выполняется быстрее, чем время создания потока (единицы микросекунд), использование `std::thread` напрямую приведет к деградации производительности.

Memory Overhead: Стек и Виртуальная память

По умолчанию в Linux размер стека для главного потока и создаваемых потоков (`pthread_create`) составляет **8 МБ**. Это значение можно проверить командой `ulimit -s`.

Важно!

Поток не потребляет 8 МБ физической памяти (RAM) мгновенно. Стек выделяется в *виртуальном адресном пространстве*. Физические страницы (frames) выделяются ядром лениво (Demand Paging) по мере обращения к адресам стека.

Однако, даже виртуальная память – ресурс исчерпаемый. Рассмотрим задачу C10K (обслуживание 10 000 соединений). Если мы используем модель "один поток на соединение" (Thread-per-connection):

$$10\,000 \text{ threads} \times 8 \text{ MB} \approx 80 \text{ GB Virtual Memory}$$

Для 32-битных систем это гарантированный крах (адресное пространство всего 4 ГБ). Для 64-битных систем это создает огромную нагрузку на TLB (Translation Lookaside Buffer) процессора и структуры ядра, управляющие памятью (VMA – Virtual Memory Areas).

Если программа под нагрузкой реально использует стек (глубокая рекурсия или большие буферы на стеке), то потребление физической памяти также станет колоссальным, что приведет к ООМ (Out Of Memory) Killer.

Планировщик Linux (CFS)

В современных ОС используется **вытесняющая многозадачность** (preemptive multitasking). Разработчик не контролирует, когда поток начнет исполняться и когда будет прерван. Этим занимается планировщик ядра.

В Linux используется **CFS (Completely Fair Scheduler)**.

Механика работы CFS

CFS моделирует "идеальный многозадачный процессор" на реальном железе.

- Каждому потоку присваивается `vruntime` (virtual runtime) – время, которое поток уже провел на процессоре, взвешенное на его приоритет (nice value).
- Потоки организованы в **красно-черное дерево** (Red-Black Tree), отсортированное по `vruntime`.
- Планировщик всегда выбирает поток с наименьшим `vruntime` (самый левый узел дерева).

Квант времени и Переключение контекста

Минимальный интервал перепланирования (latency target) в Linux может составлять, например, 6-24 мс, но реальные кванты времени динамические.

Переключение контекста (Context Switch) – это не только сохранение регистров. Это:

1. Смена режима процессора (User mode → Kernel mode).
2. Загрязнение кэша инструкций и данных (Cache Pollution). Когда новый поток начинает исполнение, нужные ему данные, скорее всего, отсутствуют в L1/L2 кэшах.
3. Накладные расходы на обновление структур планировщика.

На заметку

Типичный сервер может иметь 256 логических ядер (например, AMD EPYC). Балансировка нагрузки и миграция потоков между NUMA-нодами в таких условиях – сложнейшая задача. Миграция потока на другое ядро аннулирует его прогретый кэш.

Антипаттерн: Thread per Request

Исторически веб-серверы писались по модели: "Пришел запрос → Создали поток → Обработали → Уничтожили поток".

Почему это плохо в современном HighLoad:

1. **Unbounded Concurrency:** При всплеске трафика (DDoS или хабраэффект) создается бесконечное число потоков.
2. **Thrashing:** Система тратит больше времени на переключение контекста между тысячами потоков, чем на полезную работу.
3. **Latency Tail:** Время ответа начинает непредсказуемо расти из-за очередей в планировщике.

Thread Pool (Пул потоков)

Архитектурный паттерн, при котором создается фиксированное число долгоживущих потоков (обычно равное числу аппаратных ядер), разбирающих задачи из общей очереди.

Это унифицирует место создания потоков и ограничивает их максимальное число, предотвращая перегрузку системы.

Проблемы стандартных абстракций C++

Стандарт C++ предлагает высокоуровневые инструменты, которые часто вводят в заблуждение своей простотой.

std::execution::par (C++17)

Алгоритмы с политиками исполнения, например `std::for_each(std::execution::par, ...)`, обещают автоматическое распараллеливание.

```

1 std::vector<int> data = ...;
2 std::for_each(std::execution::par, data.begin(), data.end(), [](int& x) {
3     heavy_computation(x);
4 });

```

Проблема: Стандарт не специфицирует, как именно реализуется параллелизм.

- Реализация может создать новый поток на каждый вызов алгоритма.

- Если вы вызовете такой алгоритм внутри уже запущенного параллельного кода, произойдет комбинаторный взрыв числа потоков (Oversubscription).
- Отсутствует контроль над тем, в каком пуле исполняются задачи.

В продакшен-коде использование `std::execution::par` без глубокого понимания реализации стандартной библиотеки конкретного компилятора (libstdc++, libc++) считается небезопасным.

std::async и std::future

`std::async` может запускать задачу синхронно (в том же потоке) или асинхронно, в зависимости от флагов и реализации. Возвращаемый `std::future` в деструкторе блокирует поток до завершения задачи, что часто приводит к неожиданным последовательным исполнениям там, где ожидалась параллельность.

Резюме раздела

- Создание потока – дорогая операция (>5 мкс). Избегайте частого создания/уничтожения потоков.
- Каждый поток резервирует 8 МБ виртуального адресного пространства. 100k потоков недопустимы.
- Переключение контекста стоит дорого из-за cache misses.
- Используйте Thread Pool для контроля конкурентности.
- Осторожно с `std::async` и `std::execution::par` – они скрывают детали управления потоками, что опасно для HighLoad.

Глава 2

Building a Production-Grade ThreadPool

Осознав стоимость создания потоков и ограничения операционной системы в предыдущей главе, мы приходим к необходимости архитектурного паттерна **ThreadPool**. Идея проста: мы создаем фиксированный набор потоков (worker threads) на старте приложения и переиспользуем их для выполнения множества задач.

Однако реализация надежного, готового к продакшну пула потоков — задача нетривиальная. Она требует глубокого понимания примитивов синхронизации, управления жизненным циклом потоков и обработки исключительных ситуаций. В этой главе мы шаг за шагом спроектируем и реализуем ThreadPool на современном C++.

Фундамент: Блокирующая очередь (Blocking Queue)

Сердцем любого пула потоков является очередь задач. Воркеры (потребители) должны извлекать задачи из очереди, а клиенты (производители) — добавлять их.

Ключевое требование к такой очереди — она должна быть **блокирующей** для потребителей. Если задач нет, поток-воркер не должен крутиться в цикле (`while(empty)`), потребляя 100% CPU (busy wait). Он должен "уснуть" на уровне ядра ОС, освободив ресурсы процессора, и проснуться только тогда, когда появится новая задача или поступит сигнал остановки.

Для реализации этого механизма нам потребуются три компонента:

1. `std::deque<T>` или `std::queue<T>` — контейнер для хранения данных.
2. `std::mutex` — для обеспечения взаимного исключения (mutual exclusion) при доступе к контейнеру.
3. `std::condition_variable` — для организации ожидания и уведомления потоков.

Реализация UnboundedBlockingQueue

Рассмотрим реализацию неограниченной (unbounded) очереди. "Неограниченная" означает, что метод Push никогда не блокируется из-за переполнения (пока есть оперативная память).

```
1 #include <mutex>
2 #include <condition_variable>
```

```

3 #include <deque>
4 #include <optional>
5
6 template <typename T>
7 class UnboundedBlockingQueue {
8 public:
9     void Push(T item) {
10         {
11             std::lock_guard<std::mutex> lock(mutex_);
12             buffer_.push_back(std::move(item));
13         }
14         // Уведомляем один спящий поток (если есть)
15         not_empty_.notify_one();
16     }
17
18     std::optional<T> Pop() {
19         std::unique_lock<std::mutex> lock(mutex_);
20
21         // Ожидаем, пока очередь не станет непустой ИЛИ не будет закрыта
22         not_empty_.wait(lock, [this]() {
23             return !buffer_.empty() || is_closed_;
24         });
25
26         if (buffer_.empty() && is_closed_) {
27             return std::nullopt; // Очередь закрыта и пуста
28         }
29
30         T item = std::move(buffer_.front());
31         buffer_.pop_front();
32         return item;
33     }
34
35     void Close() {
36         {
37             std::lock_guard<std::mutex> lock(mutex_);
38             is_closed_ = true;
39         }
40         // Будим ВСЕ потоки, чтобы они увидели флаг is_closed_
41         not_empty_.notify_all();
42     }
43
44 private:
45     std::deque<T> buffer_;
46     std::mutex mutex_;
47     std::condition_variable not_empty_;
48     bool is_closed_ = false;
49 };

```

Разбор механики синхронизации

1. Lock Guard vs Unique Lock

В методе Push используется `std::lock_guard`, так как нам нужно просто захватить мьюнекс на время скоупа. В методе Pop используется `std::unique_lock`. Это критически важно.

но, так как `std::condition_variable::wait` требует именно `unique_lock`. Внутри `wait` атомарно происходит следующее: 1. Мьютекс разблокируется. 2. Поток переводится в режим ожидания (`sleep`). При пробуждении мьютекс снова захватывается.

2. Spurious Wakeups (Ложные пробуждения)

Обратите внимание на конструкцию вызова `wait`:

```
1 not_empty_.wait(lock, [this]() {
2     return !buffer_.empty() || is_closed_;
3});
```

Это эквивалентно циклу:

```
1 while (!(!buffer_.empty() || is_closed_)) {
2     not_empty_.wait(lock);
3 }
```

Важно!

Никогда не используйте `if` с `wait`. POSIX Threads (и стандарт C++) допускают **Spurious Wakeups** – ситуация, когда поток просыпается без вызова `notify`. Если не проверить условие (предикат) повторно в цикле, поток может попытаться извлечь элемент из пустой очереди, что приведет к UB или исключению.

3. Graceful Shutdown (Закрытие очереди)

Метод `Close` устанавливает флаг `is_closed_` и вызывает `notify_all()`. Это необходимо для корректного завершения работы. Представим ситуацию: 5 воркеров спят в методе `Pop`, ожидая задач. Если мы просто перестанем посыпать задачи, воркеры никогда не проснутся и программа зависнет при выходе. Вызов `notify_all()` будит их всех, они проверяют предикат `is_closed_`, видят `true` и корректно выходят, возвращая `std::nullopt`.

Архитектура ThreadPool

Теперь, имея надежную очередь, мы можем реализовать сам пул.

Интерфейс задач

Для простоты будем считать, что задача – это `std::function<void()>`. Это стирает тип функтора (лямбда, указатель на функцию, объект с `operator()`), позволяя хранить разнородные задачи в одной очереди.

На заметку

В реальных высокопроизводительных системах `std::function` может быть заменена на кастомную реализацию `delegate` без аллокаций памяти (SBO – Small Buffer Optimization), так как `std::function` может аллоцировать память в куче.

Полная реализация ThreadPool

```
1 #include <vector>
2 #include <thread>
3 #include <functional>
4 #include <iostream>
5
6 // Используем UnboundedBlockingQueue, описанную выше
7 using Task = std::function<void()>;
8
9 class ThreadPool {
10 public:
11     explicit ThreadPool(size_t num_threads) {
12         Start(num_threads);
13     }
14
15     ~ThreadPool() {
16         Stop();
17     }
18
19     // Запрещаем копирование и перемещение для простоты
20     ThreadPool(const ThreadPool&) = delete;
21     ThreadPool& operator=(const ThreadPool&) = delete;
22
23     void Submit(Task task) {
24         queue_.Push(std::move(task));
25     }
26
27 private:
28     void Start(size_t num_threads) {
29         for (size_t i = 0; i < num_threads; ++i) {
30             workers_.emplace_back([this]() {
31                 WorkerRoutine();
32             });
33         }
34     }
35
36     void Stop() {
37         // 1. Сообщаем очереди, что больше задач не будет
38         // и будим всех спящих воркеров.
39         queue_.Close();
40
41         // 2. Ждем завершения каждого потока.
42         for (auto& worker : workers_) {
43             if (worker.joinable()) {
44                 worker.join();
45             }
46         }
47     }
48
49     void WorkerRoutine() {
50         while (true) {
51             // Блокируемся в ожидании задачи
52             auto task_opt = queue_.Pop();
```

```

53
54     // Если вернулся nullopt, значит очередь закрыта и пуста -> выход
55     if (!task_opt.has_value()) {
56         break;
57     }
58
59     auto& task = task_opt.value();
60     try {
61         task(); // Выполнение задачи
62     } catch (const std::exception& e) {
63         // Критически важно ловить исключения!
64         // Иначе исключение покинет поток и вызовет std::terminate
65         std::cerr << "Worker exception: " << e.what() << std::endl;
66     } catch (...) {
67         std::cerr << "Worker unknown exception" << std::endl;
68     }
69 }
70 }
71
72 private:
73     std::vector<std::thread> workers_;
74     UnboundedBlockingQueue<Task> queue_;
75 };

```

Анализ Corner Cases и Ошибок

1. Исключения в воркерах

В методе `WorkerRoutine` блок `try-catch` обязателен. Стандарт C++ гласит: если исключение покидает функцию верхнего уровня потока (в данном случае лямбду, переданную в конструктор `std::thread`), вызывается `std::terminate()`, который аварийно завершает **весь процесс**, а не только один поток.

Без `try-catch`, одна сбойная задача ("poison pill") убила бы всё приложение.

2. Порядок остановки (Destruction Order)

В деструкторе (или методе `Stop`) порядок действий критичен:

1. Сначала **Queue Close**. Это устанавливает флаг и будит потоки.
2. Только потом **Thread Join**.

Если поменять местами или забыть `Queue Close`, вызов `worker.join()` заблокируется навечно (deadlock), так как главный поток будет ждать завершения воркера, а воркер будет спать внутри `queue.Pop()`, ожидая задач, которые никогда не придут.

3. Невозможность принудительного убийства (Thread Cancellation)

Новички часто спрашивают: "Как убить задачу, которая зависла?". В стандартном C++ (да и в большинстве ОС) безопасного способа принудительно остановить поток извне **не существует**.

- `pthread_cancel` или `TerminateThread` (WinAPI) существуют, но их использование опасно.
- Если убить поток, пока он держит `std::mutex`, этот мьютекс останется заблокированным навсегда (`abandoned mutex`). Любой другой поток, попытавшийся его захватить, зависнет.
- Также не будут вызваны деструкторы стековых объектов (`RAll` не сработает), что приведет к утечкам ресурсов.

Единственный способ прервать задачу – кооперативная отмена, когда сама задача периодически проверяет флаг `std::atomic<bool> should_stop`.

4. Data Races при проверке состояния

Рассмотрим распространенную ошибку реализации очереди без мьютексов при проверке на пустоту:

```
1 // ОШИБКА!
2 bool IsEmpty() const {
3     return buffer_.empty(); // Чтение без мьютекса!
4 }
```

`std::deque::empty()` не является атомарной операцией. Чтение состояния контейнера параллельно с его модификацией (в методе `Push` из другого потока) – это классическая **Data Race**, приводящая к Undefined Behavior. Любой доступ к разделяемым данным должен быть защищен тем же мьютексом, что и операции записи.

Резюме раздела

Мы построили корректный `ThreadPool`, который:

- Использует фиксированное число потоков.
- Эффективно ожидает задачи без busy-wait (благодаря `condition_variable`).
- Корректно обрабатывает завершение работы (Graceful Shutdown).
- Устойчив к исключениям внутри пользовательских задач.

Этот код является минимальным базисом. В реальных системах (например, в игровых движках или веб-серверах) его усложняют, добавляя приоритеты задач, `work-stealing` (кражу задач между очередями разных потоков) и локальные очереди для уменьшения конкуренции за единый мьютекс.

Глава 3

User-Space Concurrency: Implementing Fibers

В предыдущих главах мы работали с потоками операционной системы (`std::thread`). Мы выяснили, что это тяжеловесные объекты: переключение контекста требует вмешательства ядра (`syscall`), загрязняет кэш и занимает микросекунды.

Для систем, требующих обработки сотен тысяч одновременных соединений (C100K+), модель "один поток на соединение" становится узким местом. Решением является **User-Space Concurrency** – реализация многозадачности силами самого приложения, без участия ядра ОС в планировании конкретных задач.

В этой главе мы реализуем библиотеку **Файбера** (Fibers, также известны как Green Threads или Coroutines) с нуля. Мы опустимся на самый низкий уровень – уровень ассемблера и регистров процессора.

Кооперативная многозадачность

В отличие от вытесняющей (preemptive) многозадачности, где ОС насилино прерывает потоки по таймеру, файбера используют **кооперативную** (cooperative) модель.

Файбер (Fiber)

Легковесный поток исполнения, управляемый планировщиком в пространстве пользователя (User Space). Переключение между файберами происходит только тогда, когда сам файбер явно отдает управление (вызывает `Yield`).

Модель M:N Threading: Мы запускаем M файберов поверх N физических потоков ОС (обычно N равно числу ядер). Если файбер блокируется (ждет I/O или мьютекс), физический поток не блокируется, а переключается на выполнение другого файбера.

Архитектура Файбера

Чтобы превратить функцию в "поток", который можно остановить и продолжить, нам нужно сохранить его состояние. Что определяет состояние исполнения?

1. **Стек (Stack):** Локальные переменные и цепочка вызовов. У файбера должен быть свой собственный стек, выделенный в куче.

2. **Регистры процессора (Context):** Значения переменных, находящихся в данный момент на "верхушке" вычислений.
3. **Указатель инструкций (RIP):** Место в коде, где мы остановились.

Инфраструктура: Стек

Стек – это просто непрерывный блок памяти. Стек растет "вниз" (от старших адресов к младшим) на архитектуре x86_64.

```

1  class Stack {
2  public:
3      explicit Stack(size_t size)
4          : stack_{std::make_unique<char[]>(size)}
5          , top_{stack_.get() + size} // Указатель на КОНЕЦ массива (старший адрес)
6      {
7          // Выравнивание стека по границе 16 байт (требование ABI)
8          top_ = reinterpret_cast<char*>(
9              reinterpret_cast<uintptr_t>(top_) & ~0xF
10         );
11     }
12
13     void* Top() const { return top_; }
14 private:
15     std::unique_ptr<char[]> stack_;
16     char* top_;
17 };

```

Deep Dive: Assembly & Calling Conventions

Самая сложная часть реализации – переключение контекста. C++ не имеет стандартных средств для прямой записи в регистры `RSP` (Stack Pointer) и `RIP` (Instruction Pointer). Нам придется использовать ассемблер.

Согласно **System V AMD64 ABI** (соглашение о вызовах в Linux/Unix), регистры делятся на две группы:

- **Caller-saved (Volatile):** `rax`, `rcx`, `rdx`, `rsi`, `rdi`, `r8`-'r11`. Функция может менять их как угодно. Если вызывающему коду они нужны, он сам их сохраняет перед вызовом.
- **Callee-saved (Non-volatile):** `rbx`, `rbp`, `r12`-'r15`. Вызываемая функция **обязана** сохранить их значения и восстановить перед возвратом.

Важно!

Для реализации переключения контекста нам достаточно сохранять только **Callee-saved** регистры. Потому что переключение контекста выглядит как вызов функции `SaveContext`. Компилятор C++, генерирующий код вызова этой функции, уже сохранил все нужные ему Volatile-регистры на стеке (или они ему не нужны). Наша задача – сохранить только те регистры, которые компилятор ожидает увидеть неизменными после возврата из функции.

Структура контекста:

```

1 struct Context {
2     void* rip; // Instruction Pointer (адрес возврата)
3     void* rsp; // Stack Pointer
4
5     // Callee-saved регистры
6     void* rbp; // Base Pointer
7     void* rbx;
8     void* r12;
9     void* r13;
10    void* r14;
11    void* r15;
12 };

```

Магия переключения: SaveContext и JumpContext

Мы реализуем две функции на ассемблере.

1. SaveContext

Сохраняет текущее состояние в структуру `Context` и возвращает `0`. Сигнатура: `extern "C" int SaveContext(Context* ctx);`

На заметку

В System V ABI первый аргумент функции (в данном случае указатель `ctx`) передается в регистре RDI. Возвращаемое значение – в RAX.

```

.global SaveContext
SaveContext:
    // RDI содержит указатель на Context

    // 1. Сохраняем адрес возврата (RIP).
    // Текущий RIP лежит наверху стека (как положила (инструкция call)).
    mov rax, [rsp]
    mov [rdi + 0], rax // ctx->rip = [rsp]

    // 2. Сохраняем указатель стека (RSP).
    // Нам нужен RSP для вызова call, поэтому +8 байт пропускаем (адрес возврата)
    lea rax, [rsp + 8]
    mov [rdi + 8], rax // ctx->rsp = rsp + 8

    // 3. Сохраняем остальные регистры
    mov [rdi + 16], rbp
    mov [rdi + 24], rbx
    mov [rdi + 32], r12
    mov [rdi + 40], r13
    mov [rdi + 48], r14
    mov [rdi + 56], r15

    // 4. Возвращаем 0 (ESaveContextResult::Saved)
    xor rax, rax
    ret

```

2. JumpContext

Восстанавливает состояние из структуры и передает управление. Эта функция **не возвращает управление** туда, откуда её вызвали. Она "возвращается" туда, где был сохранен контекст. Сигнатура: `extern "C" void JumpContext(Context* ctx);`

```
.global JumpContext
JumpContext:
    // RDI содержит указатель на Context

    // 1. Восстанавливаем Callee-saved регистры
    mov rbp, [rdi + 16]
    mov rbx, [rdi + 24]
    mov r12, [rdi + 32]
    mov r13, [rdi + 40]
    mov r14, [rdi + 48]
    mov r15, [rdi + 56]

    // 2. Восстанавливаем стек ! Самый( важный момент )
    mov rsp, [rdi + 8]

    // 3. Подготавливаем "возврат".
    // Мы хотим , чтобы поток продолжил выполнение так , будто
    // SaveContext только что вернуло управление , но с другим результатом .
    // Загружаем сохраненный RIP в стек , чтобы инструкция ret его забрала .

    mov rax, [rdi + 0] // Загружаем целевой RIP
    push rax           // Кладем его на стек

    // 4. Устанавливаем возвращаемое значение 1 (ESaveContextResult::Resumed)
    mov rax, 1

    // 5. Прыжок!
    ret // Снимет RIP с стека и прыгнет туда
```

Важно!

Атрибут `__attribute__((returns_twice))` для `SaveContext` критически важен. Он сообщает компилятору, что функция может вернуть управление дважды (как `setjmp`), запрещая агрессивные оптимизации (например, хранение переменных только в регистрах без сброса на стек перед вызовом).

Трамплин и Запуск файбера

Мы научились переключаться. Но как запустить новый файбер в первый раз? У него нет сохраненного контекста.

Мы должны сfabриковать контекст вручную.

1. Выделяем стек.
2. Устанавливаем `context.rsp` на вершину этого стека.
3. Устанавливаем `context.rip` на адрес функции-трамплина.

Зачем нужен трамплин? Мы не можем просто указать `rip` на пользовательскую функцию. Когда пользовательская функция завершится (`ret`), процессору нужно знать, куда вернуться. В пустом стеке нет адреса возврата → Segfault. Трамплин — это обертка, которая вызывает пользовательский код, а после его завершения вызывает `Scheduler::Terminate()`, корректно завершая жизнь файбера.

```

1 // static
2 void Fiber::Trampoline() {
3     // Получаем текущий файбер (установленный шедулером перед прыжком)
4     Fiber* fiber = GetCurrentFiber();
5     fiber->RunUserRoutine(); // Вызов пользовательской лямбды
6
7     // Если мы здесь, значит задача выполнена.
8     // Файбер нельзя просто удалить, нужно уйти в шедулер.
9     GetCurrentScheduler()->Terminate();
10 }
11
12 Fiber::Fiber(std::function<void()> routine)
13     : stack_(kStackSize), routine_(std::move(routine)) {
14
15     context_.rsp = stack_.Top();
16     context_.rip = reinterpret_cast<void*>(Trampoline);
17 }
```

Планировщик (Scheduler)

Планировщик управляет очередью готовых файберов (`Runnable`) и распределяет процессорное время. Поскольку мы пишем однопоточный планировщик (для простоты), он работает в главном потоке программы.

Intrusive List

Чтобы избежать аллокаций памяти при каждом добавлении файбера в очередь планировщика (что убило бы производительность), мы используем **Инtrузивный список**. Файбер сам является узлом списка:

```
1 class Fiber : public IntrusiveList

```

Это позволяет перемещать файбер между очередями (Ready Queue ↔ Wait Queue мьютекса) простым переписыванием указателей `prev/next` внутри объекта файбера, без `new/delete`.

Цикл планирования

```

1 void Scheduler::Run() {
2     while (!queue_.Empty()) {
3         Fiber* fiber = queue_.PopFront();
4         current_fiber_ = fiber;
5     }
}
```

```

6   // Магия переключения:
7   // 1. Сохраняем контекст планировщика в main_context_
8   // 2. Загружаем контекст файбера
9   Switch(&main_context_, fiber->GetContext());

10  // Сюда мы попадаем, когда файбер вернул управление (Yield/Suspend)

11  if (fiber->GetState() == EFiberState::Finished) {
12      delete fiber;
13  } else if (fiber->GetState() == EFiberState::Runnable) {
14      queue_.PushBack(fiber);
15  }
16  // Если Suspended, то ничего не делаем (он хранится в wait queue
17  // примитива)
18
19 }
20 }
```

Примитивы синхронизации: Yield и Suspend

В мире файберов нельзя использовать `std::mutex`, так как он заблокирует весь физический поток, остановив все файбера планировщика. Нам нужны свои примитивы.

Yield (Добровольная уступка)

Файбер говорит: "Я еще не закончил, но готов уступить место другим".

1. Состояние файбера остается `Runnable`.
2. Переключаемся в контекст шедулера.
3. Шедулер кладет файбер в конец очереди.

Suspend (Блокировка)

Используется в мьютексах и condition variables.

1. Файбер меняет состояние на `Suspended`.
2. Переключаемся в контекст шедулера.
3. **Важно:** Шедулер НЕ кладет файбер обратно в очередь `RunQueue`. Файбер "исчезает" для планировщика.
4. Ответственность за возвращение файбера к жизни берет на себя тот, кто его заблокировал (например, `Mutex::Unlock` вызовет `Schedule(fiber)`).

Резюме раздела

Мы реализовали файберы – механизм User-Space многозадачности.

- Мы контролируем стек и регистры вручную.
- Переключение контекста стоит десятки наносекунд (просто `mov` регистров), в отличие от микросекунд у потоков ОС.
- Мы используем только Callee-saved регистры, полагаясь на System V ABI.
- Кооперативная многозадачность требует явного `Yield` или ожидания на примитивах синхронизации, чтобы система была отзывчивой.

Эта технология лежит в основе Goroutines в Go, Coroutines в C++20 и Fiber API в Windows.

Глава 4

Linux Low-Level Sync: The Futex

До сих пор мы рассматривали два полюса синхронизации:

1. **Спинлоки (Spinlocks)**: Работают полностью в User Space, используя атомарные инструкции. Они феноменально быстры при отсутствии конкуренции, но при ожидании "сжигают" процессорное время впустую, нагревая воздух.
2. **Тяжелые мьютексы (Kernel Mutexes)**: Старые реализации (например, в ранних версиях Linux) всегда обращались к ядру ОС для блокировки. Это экономит CPU (поток спит), но системный вызов ('syscall') стоит сотни наносекунд, даже если конкуренции нет.

Истина, как всегда, посередине. Большую часть времени мьютекс свободен, и захватывать его нужно быстро (как спинлок). Но если он занят, поток должен честно уснуть (как системный мьютекс).

Именно эту гибридную модель реализует **Futex** (Fast Userspace Mutex) – фундаментальный примитив синхронизации в Linux, на котором построены `std::mutex`, `pthread_mutex` и механизмы синхронизации в JVM и Go runtime.

Философия Futex

Основная идея Futex звучит так: "**Не тревожь ядро, если нет конкуренции**".

Futex (Fast Userspace Mutex)

Механизм, позволяющий потокам синхронизироваться через общую область памяти в пространстве пользователя, прибегая к системным вызовам только в случае необходимости ожидания (блокировки) или пробуждения других потоков.

Технически, futex состоит из двух компонентов:

1. **Aligned Integer в User Space**: Обычная 32-битная переменная (обычно `std::atomic<int32_t>`), хранящая состояние блокировки.
2. **Очередь ожидания в Kernel Space**: Структура внутри ядра, где хранятся спящие потоки, привязанные к физическому адресу этой переменной.

Системный вызов futex()

В Linux нет отдельных сисколлов `futex_wait` или `futex_wake`. Есть один мультиплексированный вызов `sys_futex`. Стандартная библиотека C++ не предоставляет прямого доступа к нему, поэтому для реализации собственных примитивов приходится использовать функцию `syscall`.

```

1 #include <linux/futex.h>
2 #include <sys/syscall.h>
3 #include <unistd.h>
4 #include <atomic>
5
6 // Обертка над системным вызовом
7 long sys_futex(void* addr1, int op, int val1,
8                 const struct timespec* timeout, void* addr2, int val3) {
9     return syscall(SYS_futex, addr1, op, val1, timeout, addr2, val3);
10 }
```

Рассмотрим две главные операции.

FUTEX_WAIT

Сигнатура: `futex(addr, FUTEX_WAIT, expected_val, ...)`

Семантика: "Проверь атомарно, что по адресу `addr` все еще лежит значение `expected_val`. Если это так – усыпь текущий поток. Если значение изменилось – верни управление немедленно (код ошибки `EAGAIN`)."

Эта атомарная проверка критически важна. Рассмотрим сценарий гонки без нее (Lost Wakeup Problem):

1. Поток А видит, что мьютекс занят (`val == 1`).
2. Поток А решает уснуть.
3. В этот момент происходит переключение контекста.
4. Поток Б освобождает мьютекс (`val = 0`) и вызывает `WAKE`. Но никто еще не спит, уведомление уходит в пустоту.
5. Управление возвращается к А.
6. Поток А вызывает "безусловный sleep" и засыпает навсегда, хотя мьютекс уже свободен.

Futex решает это, выполняя проверку значения и засыпание как единую транзакцию внутри ядра (под спинлоком планировщика).

FUTEX_WAKE

Сигнатура: `futex(addr, FUTEX_WAKE, count, ...)`

Семантика: "Разбуди `count` потоков, ожидающих по этому адресу". Обычно используют:

- `count = 1`: аналог `notify_one` (для мьютексов).
- `count = INT_MAX`: аналог `notify_all` (для cond_var или `barrier`).

Реализация Mutex на базе Futex

Напишем простейший мьютекс. Мы будем использовать три состояния переменной `state`:

- 0: Разблокировано (Unlocked).
- 1: Заблокировано, нет ожидающих (Locked, no waiters).
- 2: Заблокировано, есть ожидающие (Locked, with waiters).

Состояние "2" необходимо, чтобы оптимизировать `Unlock`. Если при разблокировке мы видим "1", значит, будить никого не надо, и можно не делать дорогой сисколл `FUTEX_WAKE`.

```

1  class FutexMutex {
2  public:
3      void lock() {
4          int c;
5          // 1. Fast Path: Попытка атомарно сменить 0 -> 1.
6          // Если успешно, мы захватили мьютекс без системных вызовов.
7          if ((c = cmpxchg(0, 1)) != 0) {
8
9              // 2. Slow Path: Мьютекс занят.
10             // Если состояние было 1, меняем его на 2 (сигнализируя о наличии
11             // → ждущего).
12             if (c != 2) {
13                 c = xchg(2);
14             }
15
16             // Крутимся в цикле, пока не захватим
17             while (c != 0) {
18                 // Пытаемся уснуть.
19                 // Ядро усыпит нас ТОЛЬКО если state == 2.
20                 // Если state изменился (кто-то сделал unlock), мы не уснем.
21                 sys_futex(&state, FUTEX_WAIT, 2, nullptr, nullptr, 0);
22
23                 // Проснулись (или не спали). Пробуем захватить 2 -> 2
24                 // (фактически просто check, но нам нужно атомарно убедиться)
25                 // В реальной реализации тут стоит повторить попытку 0 -> 2
26                 c = xchg(2);
27             }
28         }
29     }
30
31     void unlock() {
32         // 1. Fast Path: Атомарно меняем любое значение на 0.
33         // fetch_sub(1) было бы оптимизацией, чтобы отличить 1 от 2,
34         // но здесь для простоты используем exchange.
35         if (state.exchange(0) == 2) {
36             // 2. Slow Path: Будим одного ждущего, только если старое значение
37             // → было 2.
38             sys_futex(&state, FUTEX_WAKE, 1, nullptr, nullptr, 0);
39         }
40     }
41     // Вспомогательные обертки над std::atomic

```

```

42     int cmpxchg(int expected, int desired) {
43         int expected_copy = expected;
44         state.compare_exchange_strong(expected_copy, desired);
45         return expected_copy;
46     }
47
48     int xchg(int desired) {
49         return state.exchange(desired);
50     }
51
52     std::atomic<int> state{0}; // Требует выравнивания (обычно 4 байта)
53 };

```

Внутри ядра: Kernel Wait Queues

Что происходит, когда вы зовете `FUTEX_WAIT`?

1. **Hash Table Lookup:** Ядро не может просто взять виртуальный адрес `&state`, так как у разных процессов разные адресные пространства (а futex может быть shared между процессами). Ядро транслирует виртуальный адрес в физический (или в уникальный ключ inode+offset для memory-mapped файлов). 2. По этому ключу вычисляется хэш, и выбирается соответствующая цепочка (bucket) в глобальной хэш-таблице `futex_queues`. 3. **Spinlock:** Блокируется бакет хэш-таблицы. 4. **Atomic Check:** Ядро читает значение из user-space памяти. Если оно не совпадает с `expected`, блокировка снимается, и функция возвращает ошибку. 5. **Sleep:** Если значения совпадают, создается объект ожидания, добавляется в очередь, и поток переводится в состояние `TASK_INTERRUPTIBLE`. Спинлок отпускается, вызывается планировщик (`schedule()`).

Этот механизм гарантирует, что futex является самым эффективным способом блокировки в Linux, сочетая скорость атомиков (в 99% случаев) и корректность системного ожидания при высокой нагрузке.

На заметку

Понимание futex необходимо для следующей главы, так как многие Lock-Free алгоритмы используют его как "fallback mechanism" (механизм отката). Нельзя бесконечно крутиться в CAS-цикле, если ожидание затягивается; рано или поздно поток нужно усыпить.

Глава 5

Lock-Free Programming: Atomic Foundations

В предыдущих главах мы использовали блокировки (мьютексы, футексы), чтобы защитить общие данные. Это подход **пессимистичной** синхронизации: "Я предполагаю, что кто-то помешает мне, поэтому я закрою дверь на замок, прежде чем что-то делать".

Lock-Free программирование – это подход **оптимистичной** синхронизации. Потоки не блокируют друг друга. Вместо этого они пытаются выполнить операцию, и если обнаруживают конфликт (кто-то другой изменил данные быстрее), они повторяют попытку.

Определения и Гарантии

Часто термин "Lock-Free" используют неправильно, называя так любой код без `std::mutex`. Формальное определение основано на гарантиях прогресса системы.

Lock-Free (Свобода от блокировок)

Алгоритм называется Lock-Free, если гарантируется, что хотя бы один поток в системе будет продвигаться вперед (делать полезную работу) за конечное число шагов, даже если другие потоки замедлились или остановились.

Wait-Free (Свобода от ожидания)

Более сильная гарантированность. Каждый поток гарантированно завершает свою операцию за конечное число шагов, независимо от действий других потоков.

Важно!

Lock-Free не означает "быстрее". Lock-Free алгоритмы часто потребляют больше CPU из-за циклов повторных попыток (CAS-loops) и могут быть медленнее качественных мьютексов при высокой конкуренции (contention). Их главная цель – устойчивость к остановке потоков (например, если поток убит или прерван планировщиком, он не держит блокировку, мешая другим).

Атомики в C++

Стандарт C++11 ввел модель памяти и тип `std::atomic<T>`.

Операции над атомиками неделимы. Невозможно увидеть атомарную переменную в "частично измененном" состоянии. Основные операции:

- `load()`: Атомарное чтение.
- `store()`: Атомарная запись.
- `exchange(val)`: Записывает новое значение и возвращает старое (RMW – Read-Modify-Write).
- `fetch_add(val)`: Прибавляет значение и возвращает старое.

Compare-And-Swap (CAS)

Фундаментом всех Lock-Free структур данных является операция **Compare-And-Swap**. В C++ она представлена двумя методами: `compare_exchange_weak` и `compare_exchange_strong`.

Сигнатура:

```
1 bool compare_exchange_strong(T& expected, T desired);
```

Логика (псевдокод, выполняемый атомарно):

```
1 if (*this == expected) {
2     *this = desired;
3     return true;
4 } else {
5     expected = *this; // Обновляем expected актуальным значением!
6     return false;
7 }
```

Паттерн CAS-Loop

Рассмотрим простейшую задачу: атомарное изменение переменной по произвольной формуле $X_{new} = f(X_{old})$. Почему нельзя просто написать `atomic = f(atomic.load())`? Потому что между чтением (`load`) и записью (`store`) может вклинииться другой поток и изменить значение. Наш `store` перезатрет его результат. Это классическая гонка (Lost Update).

Правильный подход – использование цикла CAS:

```
1 std::atomic<int> value{0};
2
3 void atomic_update(int operand) {
4     int old_val = value.load(); // 1. Снэпшот состояния
5     int new_val;
6     do {
7         // 2. Вычисление нового значения на основе локальной копии
8         new_val = old_val + operand; // Здесь может быть любая функция f(old_val)
9
10        // 3. Попытка публикации
11        // Если value все еще равно old_val, то записать new_val.
```

```

12     // Если нет (кто-то успел изменить value), то обновить old_val и вернуть
13     // → false.
14 } while (!value.compare_exchange_weak(old_val, new_val));

```

Strong vs Weak CAS

- compare_exchange_strong: Гарантирует успех, если значение равно ожидаемому. Обычно реализуется инструкцией процессора (например, `LOCK CMPXCHG` на x86).
- compare_exchange_weak: Разрешает **Spurious Failures** (ложные отказы). Может вернуть `false`, даже если значение равно ожидаемому.

Зачем нужен weak? На некоторых архитектурах (ARM, PowerPC) CAS реализуется через пару инструкций `LL/SC` (Load-Linked / Store-Conditional). Если между ними произошло прерывание или обращение к кэш-линии, `SC` может не сработать. `weak` версия позволяет избежать лишних циклов внутри самой инструкции CAS, перекладывая ответственность на внешний цикл пользователя. В циклах (`while`) всегда предпочтительнее использовать `weak`, так как цикл все равно перезапустится.

Пример: Lock-Free Stack (Treiber Stack)

Реализуем классический Lock-Free стек (LIFO). Это одна из простейших структур, так как все изменения происходят только на верхушке (`head`).

Узлы стека:

```

1 template <typename T>
2 struct Node {
3     T data;
4     Node* next;
5
6     Node(const T& d) : data(d), next(nullptr) {}
7 };

```

Операция Push

```

1 std::atomic<Node<T>*> head{nullptr};
2
3 void Push(const T& data) {
4     Node<T>* new_node = new Node<T>(data);
5
6     // Фаза 1: Читаем текущую голову
7     new_node->next = head.load(std::memory_order_relaxed);
8
9     // Фаза 2: Пытаемся подменить голову на нашу новую ноду
10    // Если head изменился (кто-то успел сделать push/pop),
11    // CAS обновит new_node->next актуальным значением head
12    // и мы попробуем снова.

```

```

13     while (!head.compare_exchange_weak(new_node->next, new_node,
14                                         std::memory_order_release,
15                                         std::memory_order_relaxed)) {
16         // Тело цикла пустое, вся работа в условии
17     }
18 }
```

Это идеально работает. Новый узел невидим для других потоков, пока CAS не увенчается успехом. Как только CAS прошел, узел мгновенно становится новой головой.

Операция Pop и утечки памяти

```

1 std::optional<T> Pop() {
2     Node<T>* old_head = head.load(std::memory_order_acquire);
3
4     while (old_head && !head.compare_exchange_weak(old_head, old_head->next,
5                                                 std::memory_order_acquire,
6                                                 std::memory_order_relaxed)) {
7         // Если CAS не прошел, old_head обновился. Проверяем, не стал ли он
8         → nullptr.
9     }
10    if (!old_head) return std::nullopt; // Стек пуст
11
12    T res = old_head->data;
13
14    // ПРОБЛЕМА: Когда удалять old_head?
15    // delete old_head; // <--- ОПАСНО!
16
17    return res;
18 }
```

Проблема безопасного удаления памяти

В коде выше строка `delete old_head` закомментирована не случайно. Это главная проблема Lock-Free структур в C++.

Представим сценарий: 1. Поток А читает `head` (адрес 0x100) в локальную переменную `old_head`. 2. Поток А засыпает перед CAS. 3. Поток Б делает `Pop`, успешно меняет `head`, забирает ноду 0x100 и вызывает `delete` (возвращает память ОС). 4. Поток А просыпается. Он пытается обратиться к `old_head->next` внутри CAS или просто сравнить значение. **5. Use-After-Free:** Память по адресу 0x100 уже освобождена и, возможно, выделена под другой объект. Чтение `old_head->next` приводит к Segfault или чтению мусора.

В языках с Garbage Collector (Java, Go, C#) этой проблемы нет – GC не удалит объект, пока на него есть ссылки (включая локальную ссылку в потоке А). В C++ нам нужно реализовать собственный механизм управления памятью для Lock-Free.

Эта проблема и её решения (Hazard Pointers, RCU) настолько объемны, что им посвящена отдельная глава.

Спинлок (Spinlock) на атомиках

Для полноты картины реализуем примитивнейший мьютекс (спинлок) на `std::atomic_flag`. Это единственный тип, который гарантированно является Lock-Free на всех архитектурах (даже самых простых).

```

1  class Spinlock {
2      // ATOMIC_FLAG_INIT устанавливает флаг в состояние false (clear)
3      std::atomic_flag flag = ATOMIC_FLAG_INIT;
4
5  public:
6      void lock() {
7          // test_and_set устанавливает флаг в true и возвращает предыдущее
8          // значение.
9          // Мы крутимся, пока предыдущее значение было true (значит, было занято).
10         while (flag.test_and_set(std::memory_order_acquire)) {
11             // Оптимизация для процессора: "я в цикле ожидания"
12             // Снижает энергопотребление и позволяет SMT-потокам работать
13             // быстрее.
14 #if defined(__x86_64__) || defined(_M_X64)
15             _mm_pause();
16 #elif defined(__aarch64__)
17             asm volatile("yield");
18 #endif
19         }
20     }
21
22     void unlock() {
23         flag.clear(std::memory_order_release);
24     }
25 };

```

Резюме раздела

- **Атомики** обеспечивают неделимость операций и видимость изменений между потоками (Memory Ordering).
- **CAS (Compare-And-Swap)** – универсальный примитив для построения Lock-Free алгоритмов.
- **CAS-Loop** позволяет применять произвольные функции к разделяемым данным оптимистично.
- Основная сложность Lock-Free не в алгоритмах изменения данных, а в **управлении памятью** (когда безопасно удалить узел?).

Глава 6

Advanced Lock-Free: ABA & Memory Reclamation

В предыдущей главе мы реализовали Lock-Free стек и столкнулись с фундаментальной проблемой отсутствия Garbage Collector в C++: мы не знаем, когда безопасно вызывать `delete` для узла, исключенного из структуры данных.

Однако проблема управления памятью в Lock-Free алгоритмах глубже, чем просто утечки или Segfault. Повторное использование памяти (memory recycling) может привести к логическому разрушению структуры данных, даже если все указатели валидны. Этот феномен известен как проблема ABA.

В этой главе мы разберем механику ABA на уровне состояний памяти, изучим аппаратные методы решения (Tagged Pointers) и реализуем программный алгоритм безопасной рекламации памяти – Hazard Pointers. Также мы рассмотрим более сложную структуру – Lock-Free очередь Майкла-Скотта.

Проблема ABA

Название ABA происходит от последовательности изменения состояний: значение было *A*, стало *B*, затем снова стало *A*. Для наивного наблюдателя, использующего CAS (Compare-And-Swap), кажется, что значение не менялось, хотя мир вокруг изменился кардинально.

Анатомия катастрофы (Сценарий на Стеке)

Рассмотрим наш Lock-Free стек (LIFO). Состояние стека: Топ → А → В → С.

Поток 1 (Жертва):

- Планирует выполнить Pop.
- Читает текущую голову: `old_head = A` (адрес 0x1000).
- Читает следующий элемент: `next = A->next` (узел В, адрес 0x2000).
- Поток прерывается планировщиком ОС.

Поток 2 (Разрушитель):

- Выполняет Pop: Успешно удаляет А. Стек: Топ → В → С.
- Удаляет узел А (`delete A`). Память по адресу 0x1000 возвращается аллокатору.

- Выполняет Pop: Успешно удаляет В. Стек: Тор → С.
- Удаляет узел В (delete B). Адрес 0x2000 свободен.
- Выполняет Push: Создает новый узел. Аллокатор, оптимизируя работу, возвращает только что освобожденный адрес 0x1000.
- Записывает в новый узел (по адресу 0x1000) какие-то данные. Его next указывает на С.
- Стек: Тор → А' (0x1000) → С.

Поток 1 (Продолжение):

- Просыпается. Готовится выполнить CAS для смены головы.
- Аргументы CAS:
 - Адрес переменной: &Тор.
 - Ожидаемое значение: 0x1000 (узел А).
 - Новое значение: 0x2000 (узел В, сохраненный в переменной next).
- **CAS Проверка:** Текущий Тор равен 0x1000? Да, равен (это новый узел А').
- **CAS Успех:** Тор меняется на 0x2000 (узел В).

Результат: Голова стека указывает на адрес 0x2000. Но узел по этому адресу был удален Потоком 2! Стек сломан. Следующий Pop приведет к чтению освобожденной памяти (Use-After-Free) или, что хуже, к чтению данных другого объекта, который аллокатор разместил по адресу 0x2000.

Решение 1: Tagged Pointers (Указатели с версиями)

Проблема АВА возникает из-за того, что указатель (адрес) не уникален во времени. Адрес 0x1000 сегодня – это узел стека, завтра – текстура видеокарты, послезавтра – снова узел стека.

Чтобы сделать указатель уникальным, нужно добавить к нему **счетчик версий** (Tag). Каждое обновление указателя (CAS) должно инкрементировать этот счетчик. Тогда последовательность $A \rightarrow B \rightarrow A$ превратится в $A_1 \rightarrow B_2 \rightarrow A_3$. CAS, ожидающий A_1 , провалится, увидев A_3 .

Реализация на x86_64

В 64-битных процессорах виртуальные адреса используют только младшие 48 бит. Старшие 16 бит обычно должны быть нулями (или копией 47-го бита). Мы можем использовать эти 16 бит для хранения счетчика.

```

1 // Упрощенная концепция (требует битхаков)
2 struct TaggedPointer {
3     uint64_t raw; // 48 бит адрес + 16 бит тег
4
5     static const uint64_t TAG_MASK = 0xFFFF000000000000;
6     static const uint64_t PTR_MASK = 0x0000FFFFFFFFFFFF;
7

```

```

8     Node* get_ptr() const {
9         return reinterpret_cast<Node*>(raw & PTR_MASK);
10    }
11
12    uint16_t get_tag() const {
13        return (raw & TAG_MASK) >> 48;
14    }
15
16    // При создании нового ptr инкрементируем tag
17    static TaggedPointer make(Node* ptr, uint16_t old_tag) {
18        return { (uint64_t(ptr) & PTR_MASK) | (uint64_t(old_tag + 1) << 48) };
19    }
20};

```

В C++ для этого удобно использовать двойную ширину CAS (`cmpxchg16b` на x86_64), который оперирует парой 64-битных чисел (итого 128 бит). Стандартный `std::atomic<shared_ptr>` часто реализован именно так, но это дорого.

Решение 2: Hazard Pointers (Опасные указатели)

Tagged Pointers решают АВА, но не решают проблему преждевременного `delete`. Самым надежным чисто программным решением в C++ (без использования GC) является метод **Hazard Pointers** (HP), предложенный Майклом (Maged Michael).

Идея: Каждый поток-читатель имеет свой личный список "опасных указателей" (обычно 1-2 слота). Перед тем как работать с указателем, поток записывает его в свой HP-слот, объявляя всем остальным: "Я читаю этот объект, не удаляйте его!".

Поток-писатель, желающий удалить объект, сначала сканирует списки HP всех потоков.

- Если адрес найден в чьем-то HP – удалять нельзя. Адрес добавляется в список "на пенсию" (Retire List).
- Если адрес нигде не найден – можно безопасно вызвать delete.

Алгоритм работы с HP в Pop

Это сложнее, чем кажется, из-за гонок данных. Мы не можем просто записать указатель в HP, потому что к моменту записи он уже мог быть удален.

```

1 // Глобальный массив Hazard Pointers (по одному слоту на поток)
2 std::atomic<Node*> HP[MAX_THREADS];
3
4 std::optional<T> Pop(int thread_id) {
5     Node* old_head;
6     Node* next_node;
7
8     while (true) {
9         // 1. Читаем текущую голову
10        old_head = head.load(std::memory_order_acquire);
11        if (!old_head) return std::nullopt;
12

```

```

13    // 2. Объявляем указатель "опасным" (защищаем его)
14    HP[thread_id].store(old_head, std::memory_order_seq_cst);

15
16    // 3. КРИТИЧЕСКАЯ ПРОВЕРКА!
17    // За время между шагом 1 и 2 old_head мог быть удален.
18    // Если head изменился, значит, наш HP защищает уже "мертвый" объект.
19    // Нужно начать сначала.
20    if (head.load(std::memory_order_acquire) != old_head) {
21        HP[thread_id].store(nullptr, std::memory_order_relaxed); // Снимаем
22        ← защиту
23        continue;
24    }

25    // 4. Теперь указатель безопасно защищен. Можно читать поля.
26    next_node = old_head->next; // Безопасно!

27
28    // 5. Пытаемся удалить из стека
29    if (head.compare_exchange_weak(old_head, next_node)) {
30        break; // Успех! Мы извлекли элемент.
31    }

32    // Неудача: снимаем защиту и пробуем снова
33    HP[thread_id].store(nullptr, std::memory_order_relaxed);
34
35}

36    HP[thread_id].store(nullptr, std::memory_order_release); // Мы закончили

37
38    // 6. Отправляем old_head на удаление
39    RetireNode(old_head);

40
41    return old_head->data;
42}
43}

```

RetireNode и Scan

Функция `RetireNode` не делает `delete` сразу. Она добавляет указатель в локальный буфер потока (`thread_local vector`). Когда буфер заполняется (например, достигает размера $2 \cdot N$, где N – число потоков), запускается процедура Scan.

Scan: 1. Собирает все ненулевые указатели из глобального массива `HP` в `std::set` (или сортированный вектор). 2. Проходит по локальному буферу "пенсионеров". 3. Если указатель из буфера ЕСТЬ в множестве `HP` – оставляем его на потом. 4. Если указателя НЕТ в множестве `HP` – делаем реальный delete.

Это гарантирует, что мы никогда не удалим объект, который кто-то читает.

Решение 3: RCU (Read-Copy-Update)

Для структур, где чтений много, а записей мало, HP может быть слишком дорогим из-за необходимости писать в `atomic` (это сбивает кэш-линии). Альтернатива – RCU.

Идея: разделить время на "эпохи" (generations). 1. У нас есть глобальный счетчик эпох. 2. Писатель, удаляющий данные, "публикует" новую версию структуры данных, а старую

откладывает в список удаления текущей эпохи. 3. Писатель ждет, пока все читатели, начавшие работу в старой эпохе, завершат свои операции (Quiescent State). 4. После этого старую эпоху можно безопасно чистить.

В user-space RCU сложен в реализации (нужно знать, когда читатель "вышел" из критической секции), но библиотека `liburcu` предоставляет готовые решения.

Практика: Lock-Free Очередь (Queue)

Очередь (FIFO) сложнее стека, так как имеет две точки изменения: `Head` (откуда забираем) и `Tail` (куда добавляем). Классический алгоритм — **Michael-Scott Queue**.

Проблемы двух указателей

Если мы просто храним `atomic<Node*> Head` и `Tail`, мы не можем обновить их оба atomарно одной инструкцией. Вставка в конец требует двух шагов:

1. `Tail->next = new_node` (линковка).
2. `Tail = new_node` (продвижение хвоста).

Если поток заснет между 1 и 2, `Tail` останется указывать на предпоследний элемент. Следующий поток, пришедший делать `Push`, увидит, что `Tail->next != nullptr`.

Идея помощи (Helping): Lock-Free алгоритм очереди устроен так, что если поток видит "оставшийся" `Tail` (у которого `next` не null), он не ждет, а **помогает** завершить операцию, продвигая `Tail` вперед с помощью CAS, и только потом делает свою вставку.

Dummy Node (Фиктивный узел)

Чтобы избежать гонок на пустой очереди (когда `Head == Tail == nullptr`), очередь Майкл-Скотта всегда содержит хотя бы один "фиктивный" узел. Изначально: `Head = Tail = new Node(dummy)`.

- Push: Добавляет после `Tail`, двигает `Tail`.
- Pop: Читает `Head->next`. Если есть, возвращает данные, и делает `Head = Head->next`. Старый `Head` (бывший dummy) удаляется (через HP), а узел, который содержал данные, становится новым dummy.

Резюме раздела

- Проблема АВА разрушает логику Lock-Free алгоритмов при переиспользовании памяти.
- **Tagged Pointers** решают АВА аппаратно, добавляя счетчик версий в неиспользуемые биты указателя.
- **Hazard Pointers** — де-факто стандарт для управления памятью в C++ Lock-Free. Они гарантируют безопасность доступа, но требуют накладных расходов на `store-load` барьера при каждом чтении.
- Реализация очередей требует кооперации потоков: если один "застрял" посередине операции, другие должны ему помочь, чтобы сохранить свойство Lock-Free.