

Курс: *Архитектура компьютера и ОС*
Лекция 1: Введение в ОС и системные вызовы

Лектор: Олег

Оглавление

1	1 Лекция	2
1.1	Введение и организационные моменты	3
1.1.1	Формула оценки	3
1.1.2	Работа с домашними заданиями	3
1.2	Зачем нужна операционная система?	3
1.2.1	Проблема прямого доступа к оборудованию	4
1.2.2	Решение: операционная система как абстракция	4
1.3	Работа с памятью в C++: краткое повторение	6
1.3.1	Линейно адресуемая память и указатели	6
1.3.2	Динамическая память	6
1.3.3	Разделение аллокации и конструирования	6
1.3.4	Арифметика указателей	7
1.4	Взаимодействие с ОС: системные вызовы	7
1.4.1	Системный вызов <code>read</code>	7
1.4.2	Системный вызов <code>write</code>	8
1.4.3	Обработка ошибок и частичных операций	10
1.5	Работа с файлами	10
1.5.1	Системные вызовы <code>open</code> и <code>close</code>	11
1.5.2	Пример чтения из файла	11
2	2 Лекция	13
2.1	Взаимодействие с носителями информации	14
2.1.1	Почему не работать с диском напрямую?	14
2.2	Права доступа в Linux	15
2.2.1	Чтение вывода <code>ls -l</code>	15
2.2.2	Пользователь, группа и остальные	15
2.2.3	Команда <code>chmod</code>	15
2.3	Файловые дескрипторы и системные вызовы	16
2.3.1	Системный вызов <code>open</code>	17
2.3.2	Структура открытого файла в ядре	18
2.3.3	Другие важные системные вызовы	18
2.4	Практика: Перенаправление ввода-вывода	20
2.5	Работа с директориями	21

3	3 Лекция	24
3.1	Дополнительные инструменты для работы с файловой системой	25
3.1.1	Новые флаги для системного вызова <code>open</code>	25
3.1.2	Получение метаданных о файлах: семейство <code>stat</code>	25
3.2	Управление памятью: виртуальная адресация	26
3.2.1	Проблема модели линейной памяти	26
3.2.2	Виртуальная и физическая память	27
3.2.3	Страничная организация памяти	27
3.3	Системные вызовы для управления памятью: <code>mmap</code>	28
3.3.1	Аргументы и флаги <code>mmap</code>	28
3.3.2	Примеры использования <code>mmap</code>	29
3.3.3	Освобождение памяти: <code>munmap</code>	30
3.4	Аргументы командной строки и переменные окружения	31
3.4.1	Аргументы командной строки	31
3.4.2	Переменные окружения	31
4	4 Лекция	33
4.1	Углублённая работа с памятью	34
4.1.1	Механизм Page Fault и ленивое выделение памяти	34
4.1.2	Структура таблиц страниц (Page Tables)	35
4.1.3	Дополнительные системные вызовы для работы с памятью	35
4.1.4	Swap (своп) и его проблемы	36
4.2	Управление процессами	38
4.2.1	Подмена процесса: семейство <code>exec</code>	38
4.2.2	Создание процесса: <code>fork</code>	38
4.2.3	Жизненный цикл процесса	39
4.2.4	Паттерн <code>fork-exec</code>	40
4.3	Межпроцессное взаимодействие: Pipelines	41
5	5 Лекция	43
5.1	Управление процессами: группы и сигналы	44
5.1.1	Группы процессов (Process Groups)	44
5.1.2	Сигналы (Signals)	44
5.1.3	Отправка сигналов и ожидание процессов	45
5.1.4	Управление памятью при <code>fork()</code>	46
5.2	Представление данных	47
5.2.1	Текстовые и бинарные форматы	47
5.2.2	Кодировки текста: от ASCII до Unicode	47
5.2.3	Unicode и его кодировки	48
5.2.4	Работа с Unicode в C/C++	49
6	6 Лекция	52

6.1	Представление целых чисел	53
6.1.1	Беззнаковые числа	53
6.1.2	Знаковые числа: Прямой код (Sign-Magnitude)	53
6.1.3	Знаковые числа: Дополняющий код (Two's Complement)	53
6.2	Выравнивание данных в памяти	55
6.2.1	Зачем нужно выравнивание?	55
6.3	Процесс сборки программы	55
6.3.1	Препроцессор и <code>#include</code>	55
6.3.2	Единицы трансляции и ускорение сборки	56
6.3.3	Объектные файлы и символы	57
6.3.4	Линковка и релокации	57
6.3.5	Формат ELF	57
6.4	Особенности C++: Имена и переменные	58
6.4.1	Искажение имён (Name Mangling)	58
6.4.2	Extern C	58
6.4.3	Глобальные переменные: <code>extern</code> против <code>static</code>	58
6.5	Основы ассемблера и архитектуры	59
6.5.1	Архитектура фон Неймана	59
6.5.2	Регистры и память	59
6.5.3	Стек и вызовы функций	59
6.5.4	Соглашение о вызовах (ABI)	59
6.6	Практика: написание функций на ассемблере	60
6.6.1	Пример 1: Возврат константы	60
6.6.2	Пример 2: Identity (аргумент -> возврат)	60
6.6.3	Пример 3: Сложение (два аргумента)	61
6.6.4	Пример 4: Условный переход (If/Else)	61
6.6.5	Пример 5: Цикл (Sum)	61
7	7 Лекция	63
7.1	Адресация памяти в x86-64	64
7.1.1	Синтаксис Scale-Index-Base (SIB)	64
7.1.2	Указание размера операнда	64
7.2	Инструкция LEA (Load Effective Address)	65
7.2.1	LEA как оптимизация компилятора	65
7.3	Работа со стеком и локальными переменными	66
7.3.1	Проблема: Callee-clobbered регистры	66
7.3.2	Решение 1: Сохранение на стеке	66
7.3.3	Решение 2: Callee-saved регистры	67
7.4	Фреймовые указатели (Frame Pointers)	67
7.4.1	Структура стека с RBP	68
7.5	Секции данных в ассемблере	69

7.5.1	Директивы ассемблера для данных	69
7.6	Флаги процессора и условные переходы	70
7.7	Взаимодействие ассемблера и C/C++	71
7.7.1	Позиционно-независимый код (PIC) и RIP-адресация	71
7.7.2	Оптимизация хвостового вызова (TCO)	71
7.7.3	Вызов функций C (scanf / printf)	72
7.8	Синтаксисы ассемблера: Intel vs. AT&T	74
8	8 Лекция	75
8.1	Оптимизация ассемблерного кода	76
8.1.1	Проблема раздельной компиляции	76
8.1.2	Встроенный ассемблер (GNU Inline Assembly)	76
8.1.3	Оптимизация на этапе компоновки (LTO)	77
8.2	Взаимодействие с ядром и побочные эффекты	79
8.2.1	Прямой вызов <code>syscall</code>	79
8.2.2	<code>volatile</code> и побочные эффекты	80
8.2.3	Использование <code>asm</code> для барьеров компиляции	80
8.3	Указатели, функции и полиморфизм	82
8.3.1	Указатели на функции и косвенные переходы	82
8.3.2	Защита от атак: <code>endbr64</code>	82
8.3.3	Реализация виртуальных функций C++	83
8.3.4	JIT-компиляция (Just-in-Time)	84
8.4	Динамическая компоновка	86
8.4.1	Мотивация и основы (.so)	86
8.4.2	Механизмы PLT и GOT	86
8.4.3	Перехват вызовов (LD_PRELOAD)	87
8.4.4	Ручная загрузка библиотек (dlopen)	88
8.5	Freestanding: Программы без <code>stdlib</code>	90
8.5.1	<code>hosted</code> vs <code>freestanding</code>	90
8.5.2	Точка входа <code>_start</code>	90
8.5.3	Загрузка и расширение знака	91
8.6	Введение в архитектуру процессора	93
8.6.1	Проблема доступа к памяти и кэши	93
8.6.2	Конвейер инструкций (Pipeline)	93
9	9 Лекция	96
9.1	Оптимизации в современных процессорах	97
9.1.1	Ассоциативность кэша и её влияние на производительность	97
9.1.2	Конвейерное и внеочередное исполнение	98
9.1.3	Спекулятивное исполнение и уязвимости	98
9.1.4	Предсказание ветвлений (Branch Prediction)	98
9.2	Представление нецелых чисел	100

9.2.1	Числа с фиксированной точкой (Fixed-Point)	100
9.2.2	Стандарт IEEE 754: числа с плавающей запятой	100
9.2.3	Специальные случаи	100
9.2.4	Погрешности и работа в C++	101
9.3	Основы многопоточности	103
9.3.1	Процессы и потоки	103
9.3.2	Синхронизация и доступ к общей памяти	103
9.3.3	Атомарные операции (<code>std::atomic</code>)	103
9.3.4	Примитивы блокирующей синхронизации	104
9.4	Классическая проблема: обедающие философы	106
9.4.1	Постановка задачи	106
9.4.2	Взаимоблокировка (Deadlock)	106
10	10 Лекция	107
10.1	Введение: Сигналы и наблюдаемый параллелизм	108
10.1.1	Наблюдаемый параллелизм и квантование	108
10.2	Аппаратная поддержка многопоточности	108
10.2.1	Hyper-threading	108
10.2.2	Привязка к ядрам (CPU Affinity)	108
10.2.3	Закон Амдала	109
10.3	Thread Safety (Потокобезопасность)	109
10.3.1	Контракты стандартной библиотеки C++	109
10.3.2	Свободные функции	110
10.4	Примитивы синхронизации	110
10.4.1	Семафор (Semaphore)	110
10.4.2	RWLock (Read-Write Lock)	110
10.4.3	Барьер (Barrier)	110
10.5	Thread Local Storage (TLS)	111
10.6	Shared Pointer и безопасность	111
10.7	Атомики и модель памяти	111
10.7.1	Compare-And-Swap (CAS)	111
10.8	Кэши и когерентность	112
10.8.1	Иерархия памяти и протокол MESI	112
10.8.2	False Sharing (Ложное разделение)	112
10.9	Процессы и Fork в многопоточной среде	113
10.10	Futex (Fast Userspace Mutex)	113
11	11 Лекция	115
11.1	Введение	116
11.2	Физический и Канальный уровни	116
11.2.1	Принципы передачи сигнала	116
11.2.2	Протокол Ethernet и MAC-адресация	116

11.3	Сетевой уровень (Network Layer)	117
11.3.1	Протокол IP и маршрутизация	117
11.3.2	Подсети и маски (CIDR)	117
11.4	Транспортный уровень (Transport Layer)	118
11.4.1	Сравнение TCP и UDP	118
11.4.2	Порядок байт (Endianness)	118
11.5	Прикладной уровень (Application Layer)	119
11.5.1	HTTP и DNS	119
11.6	Socket API и Модели конкурентности	119
11.6.1	Базовый жизненный цикл TCP-сервера	119

Глава 1

1 Лекция

1.1 Введение и организационные моменты

Этот курс посвящён изучению пользовательской части **Операционная система (ОС)** и архитектуры компьютера. Основная цель — понять, как программы выполняют свои действия на низком уровне, «под капотом» стандартных библиотечных функций. Курс рассчитан на один семестр и является обязательным, с возможностью выбрать продолжение во втором семестре.

1.1.1 Формула оценки

Итоговая оценка за курс формируется по следующей формуле:

$$\text{Оценка} = \min(10, 0.6 \cdot O_{\text{дз}} + 0.2 \cdot O_{\text{кр}} + 0.2 \cdot O_{\text{экз}} + 0.1 \cdot O_{\text{сем}}) \quad (1.1.1)$$

где:

- $O_{\text{дз}}$ — оценка за домашние задания.
- $O_{\text{кр}}$ — оценка за контрольные работы.
- $O_{\text{экз}}$ — оценка за экзамен.
- $O_{\text{сем}}$ — оценка за работу на семинарах.

Примечание

Сумма весовых коэффициентов в формуле (1.1.1) равна 1.1. Это означает, что с учётом бонусов за домашние задания можно набрать более 10 баллов, но итоговая оценка ограничивается 10 баллами.

1.1.2 Работа с домашними заданиями

Домашние задания будут выдаваться примерно раз в неделю со сроком выполнения 1–2 недели. Дедлайны «мягкие»: баллы за задание начинают убывать постепенно и достигают 10–20% от первоначальной стоимости через 3 недели после выдачи. Это сделано для того, чтобы студенты не жертвовали сном ради сдачи заданий в последний момент. Все задания будут выполняться в среде GEDLab.

1.2 Зачем нужна операционная система?

Рассмотрим простейшую программу на C++, которая считывает два числа и выводит их сумму (листинг 1.1).

```
1 #include <iostream>
2
3 int main() {
4     int a, b;
5     std::cin >> a >> b;
6     std::cout << a + b;
7     return 0;
8 }
```

Листинг 1.1 – Программа для сложения двух чисел

На первый взгляд, все операции ввода-вывода выполняются благодаря библиотеке `iostream`. Однако в самом языке C++ нет встроенных механизмов для прямого взаимодействия с устройствами, такими как экран или клавиатура. Как же тогда текст появляется на мониторе?

1.2.1 Проблема прямого доступа к оборудованию

Представим модель, в которой программа напрямую взаимодействует с аппаратными компонентами компьютера: **Центральный процессор (CPU)**, **Оперативная память (RAM)**, жестким диском, сетевой картой и т.д. (рис. 1.1).

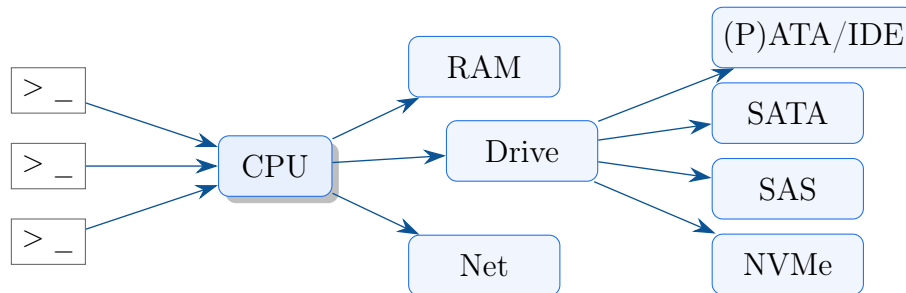


Рис. 1.1 – Модель прямого взаимодействия программы с оборудованием

Такая модель порождает две ключевые проблемы:

1. **Сложность и непереносимость.** Существует множество протоколов для взаимодействия с одним и тем же типом устройств. Например, для работы с дисками программа должна была бы поддерживать интерфейсы PATA, SATA, SAS, NVMe и другие. Аналогично, каждый производитель сетевых карт может предлагать свой уникальный протокол. Чтобы программа работала на разных компьютерах, ей пришлось бы реализовывать поддержку всех этих интерфейсов, что практически невозможно.
2. **Разделение ресурсов.** В современных системах одновременно запущены сотни и тысячи программ, в то время как количество ядер **CPU** ограничено единицами или десятками. Необходимо эффективно распределять процессорное время и другие ресурсы (память, доступ к дискам) между всеми программами. При прямом доступе программы к оборудованию сделать это было бы крайне затруднительно.

1.2.2 Решение: операционная система как абстракция

Для решения этих проблем была придумана **ОС**.

Определение: Операционная система

ОС — это программный слой, который выступает посредником между пользовательскими программами и аппаратным обеспечением компьютера.

ОС решает обе проблемы:

- Она **предоставляет унифицированный интерфейс** для работы с оборудованием. Программа работает не с конкретным жестким диском, а с абстракцией «файловой системы». ОС сама берёт на себя реализацию всех низкоуровневых протоколов.
- Она **управляет ресурсами**. ОС решает, какой программе и на какое время предоставить **CPU**, распределяет память, организует доступ к устройствам, предотвращая конфликты.

Эта модель показана на рис. 1.2.

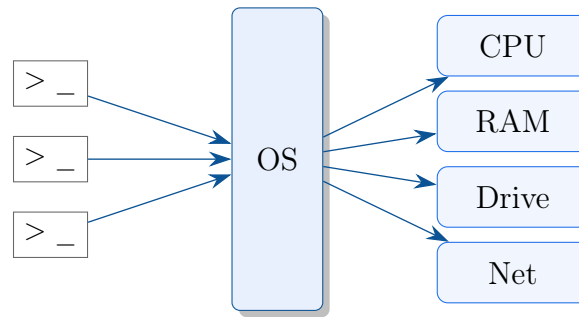


Рис. 1.2 – Операционная система как посредник

Итоги раздела

- Прямое взаимодействие программ с оборудованием сложно, непереносимо и не позволяет эффективно разделять ресурсы.
- **ОС** решает эти проблемы, предоставляя программам абстракции (файлы, сокеты) и управляя доступом к аппаратуре.

1.3 Работа с памятью в C++: краткое повторение

Взаимодействие с ОС в значительной степени происходит через память. Программа записывает данные в свою область памяти и затем просит ОС что-то с этими данными сделать. Поэтому важно освежить знания о модели памяти в C++.

1.3.1 Линейно адресуемая память и указатели

Память можно представить как большой массив байтов, где у каждого байта есть уникальный числовой адрес. Это называется моделью линейно адресуемой памяти.

Для работы с памятью используются **указатели** — переменные, которые хранят адрес. Разыменование указателя (`*ptr`) означает обращение к данным, лежащим по этому адресу. Тип указателя определяет, сколько байт будет прочитано и как они будут интерпретированы. Например, указатель типа `int64_t*` при разыменовании прочитает 8 байт и представит их как 64-битное целое число.

1.3.2 Динамическая память

Иногда память нужно выделить так, чтобы она «пережила» функцию, в которой была создана. Для этого используется динамическое выделение памяти в «куче» (heap).

```
1 // Allocating a single object
2 // Operator `new` allocates sizeof(T) bytes and constructs an object
3 T* ptr = new T{};
4
5 // Deleting an object and freeing memory
6 delete ptr;
7
8 // Allocating an array of 10 objects
9 T* arr = new T[10];
10
11 // Deleting all objects in the array and freeing memory
12 delete[] arr;
```

Листинг 1.2 – Работа с динамической памятью

Примечание

Важно соблюдать парность операторов: память, выделенную через `new`, нужно освобождать через `delete`. Память, выделенную через `new[]`, — через `delete[]`. Нарушение этого правила приводит к неопределённому поведению.

1.3.3 Разделение аллокации и конструирования

Оператор `new T` на самом деле выполняет две операции:

1. Выделение «сырой» (неинициализированной) памяти нужного размера.
2. Конструирование объекта типа `T` в этой памяти.

Эти шаги можно выполнить отдельно.

```
1 // 1. Allocate sizeof(T) raw bytes. operator new returns void*
2 void* raw_ptr = operator new(sizeof(T));
3
4 // 2. Construct an object of type T at the given address (placement new)
5 T* ptr = new (raw_ptr) T{};
```

```

6
7 // --- object `ptr` is ready to use ---
8
9 // 3. Explicitly call the destructor to destroy the object
10 ptr->~T();
11
12 // 4. Free the raw memory
13 operator delete(raw_ptr);

```

Листинг 1.3 – Явное управление памятью и объектами

Такой подход даёт больше контроля, но требует аккуратного ручного управления временем жизни объекта и памяти.

1.3.4 Арифметика указателей

В C++ арифметика указателей типизирована. Прибавление к указателю `ptr` единицы (`ptr + 1`) сдвигает его адрес не на 1 байт, а на `sizeof(*ptr)` байт, то есть к адресу следующего элемента в массиве.

- `ptr[i]` эквивалентно `*(ptr + i)`.
- Разность двух указателей одного типа `ptr1 - ptr2` даёт количество элементов (а не байт) между ними.

1.4 Взаимодействие с ОС: системные вызовы

Теперь, когда мы освежили знания о памяти, перейдём к основному механизму взаимодействия программы с ОС.

Определение: Системный вызов

Системный вызов — это основной интерфейс между пользовательскими программами и ядром ОС. Программа использует **Системный вызов**, чтобы попросить ядро выполнить действие, которое она не может выполнить сама (например, работать с файлом или сетью).

Рассмотрим два базовых системных вызова для ввода-вывода: `read` и `write`.

1.4.1 Системный вызов `read`

Функция `read` читает данные из источника, идентифицируемого **Файловый дескриптор**, в буфер.

```

1 #include <unistd.h>
2
3 ssize_t read(int fd, void *buf, size_t count);
4
5 // Example
6 char buf[10];
7 // Read from stdin (fd=0) into buf, at most 9 bytes
8 ssize_t bytes_read = read(0, buf, 9);
9
10 if (bytes_read == -1) {
11     // Error occurred
12 } else if (bytes_read == 0) {
13     // End of input (EOF)

```

```
14 } else {  
15     // Successfully read `bytes_read` bytes  
16 }
```

Листинг 1.4 – Сигнатура и использование read

Аргументы:

- `int fd`: [Файловый дескриптор](#) источника данных. По соглашению, 0 — это [Стандартный поток ввода](#).
- `void *buf`: Указатель на буфер, куда будут записаны данные.
- `size_t count`: Максимальное количество байт для чтения.

Возвращаемое значение (`ssize_t`):

- Положительное число: количество успешно прочитанных байт. **Важно:** `read` не гарантирует, что прочитает ровно `count` байт, даже если они доступны. Он может прочитать меньше.
- 0: достигнут конец файла (EOF) или потока. Больше данных для чтения нет.
- -1: произошла ошибка. Код ошибки сохраняется в глобальной переменной `errno`.

1.4.2 Системный вызов `write`

Функция `write` записывает данные из буфера в приёмник, идентифицируемый [Файловым дескриптором](#).

```
1  #include <unistd.h>  
2  
3  ssize_t write(int fd, const void *buf, size_t count);  
4  
5  // Example  
6  const char msg[] = "Hello, world!\n";  
7  // Write to stdout (fd=1), strlen(msg) bytes  
8  // We use sizeof(msg) - 1 to exclude the terminating null byte ('\0')  
9  ssize_t bytes_written = write(1, msg, sizeof(msg) - 1);  
10  
11 // Write the same message to stderr (fd=2)  
12 write(2, msg, sizeof(msg) - 1);  
13  
14 if (bytes_written == -1) {  
15     // Error occurred  
16 }
```

Листинг 1.5 – Сигнатура и использование write

Аргументы:

- `int fd`: [Файловый дескриптор](#) приёмника данных. 1 — [Стандартный поток вывода](#), 2 — [Стандартный поток ошибок](#).
- `const void *buf`: Указатель на буфер с данными для записи. Буфер константный, так как `write` его не изменяет.
- `size_t count`: Количество байт для записи.

Возвращаемое значение:

- Положительное число: количество успешно записанных байт. Как и `read`, `write` может записать меньше байт, чем было запрошено (например, если на диске закончилось место).
- -1: произошла ошибка, код которой записан в [errno](#).

1.4.3 Обработка ошибок и частичных операций

Поскольку `read` и `write` могут обработать меньше данных, чем запрошено, для надёжной передачи всего объёма данных необходимо использовать циклы.

```
1 // Writes exactly `count` bytes from `buf` to `fd`.
2 // Returns `true` on success, `false` on error.
3 bool WriteAll(int fd, const char* buf, size_t count) {
4     size_t written = 0;
5     while (written < count) {
6         ssize_t res = write(fd, buf + written, count - written);
7         if (res == -1) {
8             return false; // An error occurred
9         }
10        written += res;
11    }
12    return true;
13 }
```

Листинг 1.6 – Надёжная функция для записи всех данных

Аналогичная функция `ReadAll` должна быть реализована для чтения, но с дополнительной проверкой на возврат 0 (EOF).

Для получения текстового описания ошибки по её коду из `errno` можно использовать функцию `strerror` из заголовка `<cstring>`.

```
1 #include <cerrno>
2 #include <cstring>
3 #include <iostream>
4
5 // ... inside a function
6 if (!WriteAll(1, "Hello", 5)) {
7     // errno is set by the last failed `write` call
8     std::cerr << "Error writing data: " << strerror(errno) << std::endl;
9     return 1; // Exit with error code
10 }
```

Листинг 1.7 – Обработка ошибок с выводом сообщения

Итоги раздела

- Системные вызовы — это API операционной системы.
- `read` и `write` — базовые вызовы для неформатированного ввода-вывода.
- Файловые дескрипторы 0, 1, 2 зарезервированы для стандартных потоков `stdin`, `stdout`, `stderr`.
- Всегда проверяйте возвращаемые значения системных вызовов на ошибки (-1) и обрабатывайте частичные операции.
- Для получения информации об ошибке используйте переменную `errno`.

1.5 Работа с файлами

Стандартные потоки — это лишь частный случай. Основное применение **Файловый дескриптор** — работа с файлами на диске.

1.5.1 Системные вызовы open и close

Чтобы работать с файлом, его сначала нужно открыть с помощью системного вызова `open`.

Определение: Файловый дескриптор

Файловый дескриптор — это неотрицательное целое число, которое ОС возвращает процессу при открытии файла. Процесс использует этот дескриптор во всех последующих операциях с файлом (`read`, `write`, `close`).

```
1 #include <fcntl.h> // For flags
2 #include <unistd.h>
3
4 int open(const char *pathname, int flags, ... /* mode_t mode */);
```

Листинг 1.8 – Сигнатура системного вызова `open`

Аргументы:

- `const char *pathname`: Путь к файлу.
- `int flags`: Флаги, определяющие режим доступа (например, `O_RDONLY` — только для чтения, `O_WRONLY` — только для записи, `O_RDWR` — для чтения и записи). Флаги можно комбинировать с помощью побитового ИЛИ (`|`).
- `mode_t mode`: (Опционально) Права доступа, которые устанавливаются, если файл создаётся с флагом `O_CREAT`.

В случае успеха `open` возвращает новый **Файловый дескриптор** (обычно наименьший из доступных). В случае ошибки возвращается `-1`, а `errno` устанавливается.

После завершения работы с файлом его дескриптор необходимо освободить с помощью системного вызова `close`.

```
1 #include <unistd.h>
2
3 int close(int fd);
```

Листинг 1.9 – Сигнатура системного вызова `close`

Если не закрывать файлы, это приведёт к утечке ресурсов (файловых дескрипторов), так как их количество для одного процесса ограничено.

1.5.2 Пример чтения из файла

В листинг 1.10 показан полный цикл работы: открытие файла, чтение из него, вывод содержимого в стандартный поток и закрытие.

```
1 #include <iostream>
2 #include <fcntl.h>
3 #include <unistd.h>
4 #include <cerrno>
5 #include <cstring>
6
7 int main() {
8     const char* filename = "output.txt";
9     int fd = open(filename, O_RDONLY);
```

```
10     if (fd == -1) {
11         std::cerr << "Failed to open file " << filename << ": "
12             << strerror(errno) << std::endl;
13         return 1;
14     }
15
16     char buffer[1024];
17     ssize_t bytes_read;
18
19     // Read from file in a loop until EOF
20     while ((bytes_read = read(fd, buffer, sizeof(buffer))) > 0) {
21         // Write the read data to stdout
22         if (!WriteAll(1, buffer, bytes_read)) {
23             std::cerr << "Failed to write to stdout: "
24                 << strerror(errno) << std::endl;
25             close(fd);
26             return 1;
27         }
28     }
29
30     if (bytes_read == -1) {
31         std::cerr << "Error reading from file: " << strerror(errno) << std::endl;
32     }
33
34     close(fd); // Don't forget to close the file!
35     return 0;
36 }
37 // Assume WriteAll is defined as in listing 4.4
```

Листинг 1.10 – Чтение из файла и вывод в stdout

Итоги раздела

- Работа с файлом начинается с его открытия вызовом `open`, который возвращает [Файловый дескриптор](#).
- Полученный [Файловый дескриптор](#) используется в вызовах `read` и `write` для взаимодействия с файлом.
- После окончания работы файл необходимо закрыть вызовом `close`, чтобы освободить ресурсы.
- Каждый `open` должен иметь парный `close`, подобно паре `new/delete`.

Глава 2

2 Лекция

2.1 Взаимодействие с носителями информации

На прошлой лекции мы установили, что программы взаимодействуют с внешним миром через **Системный вызов**. Сегодня мы продолжим эту тему и углубимся во взаимодействие с **файловой системой**.

2.1.1 Почему не работать с диском напрямую?

Казалось бы, зачем нужна **файловая система**, если можно работать с жёстким диском напрямую? Тому есть две ключевые причины: сложность **Application Programming Interface (интерфейс прикладного программирования) (API)** и низкая производительность.

1. **Примитивный интерфейс.** Диск предоставляет очень аскетичное **API**: он позволяет читать и писать только «сырые» данные по указанным адресам (с такого-то по такой-то байт). В таком интерфейсе отсутствуют высокоуровневые концепции, такие как файлы, директории, права доступа и структура данных.
2. **Особенности производительности.** Жёсткий диск (HDD) — механическое устройство. Он состоит из вращающихся магнитных пластин («блинов») и считывающих головок (рис. 2.1).

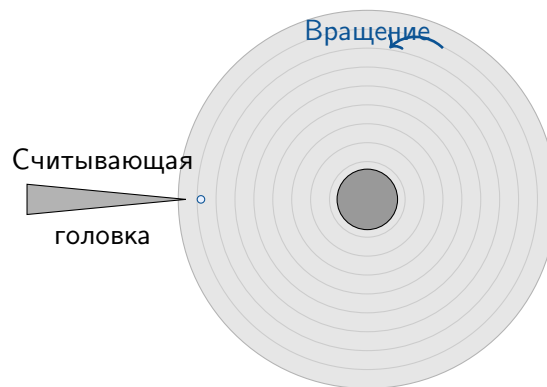


Рис. 2.1 — Упрощённая схема устройства жёсткого диска (HDD)

Скорость вращения современных дисков составляет 5000–7000 оборотов в минуту. Чтобы прочитать данные, необходимо выполнить две операции с большими задержками:

- **Позиционирование головки (seek time):** Механическое перемещение головки к нужной дорожке.
- **Ожидание вращения (rotational latency):** Ожидание, пока нужный сектор на дорожке окажется под головкой.

В среднем, ожидание нужного сектора может занимать до 5 мс. Это означает, что при чтении из случайных мест диска можно выполнить всего около 200 операций в секунду, что на порядки медленнее, чем миллиарды операций, выполняемых процессором. Для эффективной работы данные нужно располагать последовательно, минимизируя перемещения головки, но реализация такой логики — крайне сложная задача.

Определение: Файловая система

Файловая система — это уровень абстракции, предоставляемый операционной системой для организации, хранения и именования данных на носителях информа-

ции. Она скрывает сложности работы с оборудованием и предоставляет удобный и эффективный интерфейс для пользователя и программ.

2.2 Права доступа в Linux

файловая система в Linux представляет собой древовидную структуру из директорий и файлов. Для управления доступом к этим объектам используется модель прав, основанная на пользователях и группах.

2.2.1 Чтение вывода `ls -l`

Команда `ls -l` выводит подробную информацию о файлах и директориях:

```
-rw-rw-r-- 1 arch arch 4 sen 13 11:58 out
drwxr-xr-x 2 arch arch 4096 sen 13 12:00 test
```

Рассмотрим структуру вывода:

- `-rw-rw-r-`: Права доступа.
- `1`: Количество жёстких ссылок.
- `arch`: Пользователь-владелец.
- `arch`: Группа-владелец.
- `4`: Размер в байтах.
- `Сен 13 11:58`: Дата последнего изменения.
- `out`: Имя файла.

Первый символ указывает на тип: `-` для обычного файла, `d` для директории, `l` для **символическая ссылка**.

2.2.2 Пользователь, группа и остальные

Следующие 9 символов прав доступа делятся на три группы по три:

1. **Для владельца (user)**: Права пользователя, которому принадлежит файл.
2. **Для группы (group)**: Права для всех пользователей, состоящих в группе, которой принадлежит файл.
3. **Для остальных (others)**: Права для всех остальных пользователей.

Каждая тройка состоит из символов `r`, `w`, `x`:

- `r` (read): Право на чтение.
- `w` (write): Право на запись (изменение).
- `x` (execute): Право на исполнение (для программ и скриптов).

Если право отсутствует, на его месте ставится прочерк (`-`).

2.2.3 Команда `chmod`

Для изменения прав доступа используется команда `chmod` (change mode). Она поддерживает два основных синтаксиса: символический и восьмеричный.

Символический синтаксис:

```
1 # Add execute permission for the user (owner)
2 chmod u+x filename
3
4 # Remove write permission for group and others
5 chmod go-w filename
6
7 # Set permissions: read/write for user, read-only for group/others
8 chmod u=rw,go=r filename
```

Восьмеричный синтаксис: Права представляются в виде трёх восьмеричных цифр, где каждая цифра — это сумма значений для **r**, **w**, **x**:

- **r** = 4
- **w** = 2
- **x** = 1

Например, **rw-** соответствует $4 + 2 + 0 = 6$, а **r-x** — $4 + 0 + 1 = 5$.

```
1 # Corresponds to rw-rw-r-- (664)
2 chmod 664 out
3
4 # Corresponds to rwxr-xr-x (755)
5 chmod 755 script.sh
```

Примечание

Права для директорий. Права **gwx** для директорий имеют особый смысл:

- **r**: Позволяет просмотреть список файлов в директории (выполнить **ls**).
- **w**: Позволяет создавать, удалять и переименовывать файлы в директории.
- **x**: Позволяет войти в директорию (сделать **cd**) и получить доступ к файлам внутри неё (при наличии прав на сами файлы).

2.3 Файловые дескрипторы и системные вызовы

Для работы с файлами из программы операционная система предоставляет набор **Системный вызов**. Ключевой абстракцией здесь является **Файловый дескриптор**.

Определение: Файловый дескриптор

Файловый дескриптор — это неотрицательное целое число, которое процесс использует для идентификации открытого файла или другого ресурса ввода-вывода. Вместо того чтобы каждый раз передавать ядру полный путь к файлу, программа один раз вызывает **open** и получает **Файловый дескриптор**, который затем использует в вызовах **read**, **write**, **close** и др..

По умолчанию каждый процесс в Linux при запуске имеет три открытых **Файловый дескриптор**:

- 0 — стандартный поток ввода (**stdin**).
- 1 — стандартный поток вывода (**stdout**).

- 2 — стандартный поток ошибок (*stderr*).

2.3.1 Системный вызов open

Для открытия или создания файла используется [Системный вызов open](#).

```
1 #include <fcntl.h>
2 #include <sys/stat.h>
3
4 int open(const char *path, int flags, mode_t mode);
```

- path: Путь к файлу.
- flags: Битовая маска, определяющая режим доступа.
- mode: Права доступа, которые будут установлены, если файл создаётся.

Функция возвращает новый [Файловый дескриптор](#) или -1 в случае ошибки.

Основные флаги (flags):

- O_RDONLY, O_WRONLY, O_RDWR: Открыть только для чтения, только для записи или для чтения и записи. Один из этих флагов должен быть указан.
- O_CREAT: Создать файл, если он не существует.
- O_EXCL: Использовать вместе с O_CREAT. Вызов завершится ошибкой, если файл уже существует. Это позволяет атомарно создать файл и убедиться в его отсутствии до вызова.
- O_APPEND: Все операции записи будут производиться в конец файла.
- O_TRUNC: Если файл существует и открывается на запись, его содержимое усекается до нуля байт.

Создание файла и umask

При создании файла (с флагом O_CREAT) его итоговые права доступа определяются формулой:

$$\text{final_mode} = \text{mode} \& \sim \text{umask}$$

где mode — это права, переданные в open, а umask — это маска процесса. Umask определяет, какие права доступа нужно «выключить» по умолчанию. Например, если umask равна 0002 (— — w — — —), то у всех создаваемых файлов будет отбираться право на запись для «остальных».

```
1 #include <fcntl.h>
2 #include <sys/stat.h>
3 #include <unistd.h>
4
5 int main() {
6     // Create "file" if it does not exist.
7     // Error if it already exists.
8     // Permissions: read for all (0444).
9     // Mode: read and write for our process.
10    int fd = open(
11        "file",
12        O_RDWR | O_CREAT | O_EXCL,
13        S_IRUSR | S_IRGRP | S_IROTH /* 0444 */
14    );
15 }
```

```
14     );
15
16     if (fd == -1) {
17         // handle error
18         return 1;
19     }
20
21     // ... work with the file ...
22
23     close(fd);
24     return 0;
25 }
```

Листинг 2.1 – Пример использования open

Примечание

Открытый файл — это ресурс, который ядро выделяет для процесса. Как и любую другую выделенную память, его необходимо освобождать. Для этого используется [Системный вызов close\(int fd\)](#). Если этого не делать, произойдёт утечка ресурсов (файловых дескрипторов).

2.3.2 Структура открытого файла в ядре

С каждым открытым [Файловый дескриптор](#) ядро ассоциирует структуру, содержащую как минимум:

- **Флаги открытия:** Режим, в котором файл был открыт (O_RDONLY и т.д.).
- **Текущее смещение:** Позиция в файле, с которой будет происходить следующая операция чтения/записи.
- **Ссылка на inode:** Указатель на структуру файла в [файловая система](#).

Важно, что права доступа проверяются только один раз — во время вызова `open`. Все последующие операции с [Файловый дескриптор](#) (`read`, `write`) не требуют повторной проверки прав.

2.3.3 Другие важные системные вызовы

lseek: Изменение смещения

[Системный вызов lseek](#) позволяет изменить текущее [смещение](#) в файле.

```
1 #include <unistd.h>
2
3 off_t lseek(int fd, off_t offset, int whence);
```

- `fd`: [Файловый дескриптор](#), для которого меняется [смещение](#).
- `offset`: Значение смещения в байтах.
- `whence`: Точка отсчёта:
 - `SEEK_SET`: [смещение](#) отсчитывается от начала файла.
 - `SEEK_CUR`: [смещение](#) отсчитывается от текущей позиции.
 - `SEEK_END`: [смещение](#) отсчитывается от конца файла.

С помощью `lseek` можно перемещаться за конец файла. Если после такого перемещения произвести запись, то пространство между старым концом файла и новой позицией записи будет заполнено нулевыми байтами, создавая [разреженный файл](#).

`dup` и `dup2`: Копирование файловых дескрипторов

Эти вызовы создают копию [Файловый дескриптор](#).

```
1 #include <unistd.h>
2
3 int dup(int oldfd);
4 int dup2(int oldfd, int newfd);
```

`dup` создаёт копию `oldfd`, используя первый свободный номер [Файловый дескриптор](#). `dup2` создаёт копию `oldfd` с конкретным номером `newfd`. Если `newfd` уже был открыт, он атомарно закрывается.

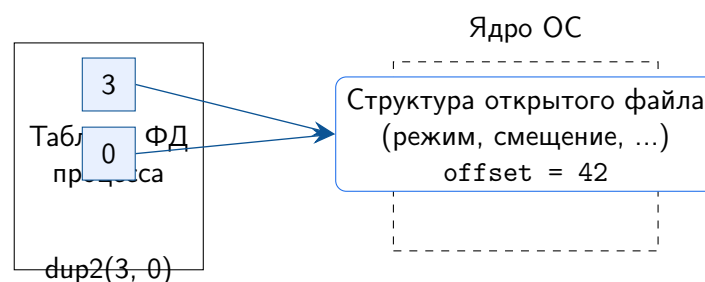


Рис. 2.2 – Схема работы `dup2`. Оба дескриптора (старый и новый) указывают на одну и ту же структуру открытого файла в ядре и разделяют общее смещение.

Ключевой особенностью является то, что новый и старый [Файловый дескриптор](#) ссылаются на одну и ту же запись в таблице открытых файлов ядра (рис. 2.2). Это означает, что они **разделяют общее смещение**: изменение позиции через один [Файловый дескриптор](#) немедленно отражается на другом.

`pipe`: Создание каналов

[Системный вызов](#) `pipe` создаёт однонаправленный [канал \(pipe\)](#) для межпроцессного взаимодействия.

```
1 #include <unistd.h>
2
3 int pipe(int pipefd[2]);
```

Вызов создаёт пару связанных [Файловый дескриптор](#) и помещает их в массив `pipefd`:

- `pipefd[0]`: [Файловый дескриптор](#) для чтения из канала.
- `pipefd[1]`: [Файловый дескриптор](#) для записи в канал.

Данные, записанные в `pipefd[1]`, можно прочитать из `pipefd[0]` в том же порядке (FIFO).

- [Файловый дескриптор](#) — это числовой идентификатор открытого ресурса.
- `open` открывает/создаёт файл и возвращает [Файловый дескриптор](#).
- `lseek` позволяет перемещаться по файлу, изменяя [смещение](#).

- `dup2` копирует [Файловый дескриптор](#), что является основой для перенаправления ввода-вывода.
- `pipe` создаёт пару [Файловый дескриптор](#) для однонаправленной передачи данных между процессами.
- Все открытые ресурсы должны быть закрыты с помощью `close`.

2.4 Практика: Перенаправление ввода-вывода

Одной из самых мощных возможностей, которую даёт `dup2`, является перенаправление стандартных потоков ввода-вывода. Рассмотрим программу, которая читает число из `stdin` и выводит его инкремент в `stdout`.

```
1 #include <iostream>
2
3 int main() {
4     int a;
5     std::cin >> a;
6     std::cout << a + 1 << std::endl;
7     return 0;
8 }
```

Листинг 2.2 – Программа с простым вводом-выводом

Мы можем перехватить её ввод и вывод, не изменяя исходный код. Для этого нужно открыть файлы для чтения и записи, а затем с помощью `dup2` подменить стандартные [Файловый дескриптор](#) (0 и 1) нашими.

```
1 #include <fcntl.h>
2 #include <unistd.h>
3 #include <string_view>
4 #include <iostream>
5
6 // Utility for error handling
7 [[noreturn]] void Fail(std::string_view msg) {
8     perror(msg.data());
9     std::abort();
10 }
11
12 void Redirect() {
13     // Open a file for reading
14     int fin = open("input.txt", O_RDONLY);
15     if (fin == -1) {
16         Fail("open input.txt");
17     }
18
19     // Open a file for writing, create if it does not exist
20     int fout = open("out.txt", O_WRONLY | O_CREAT, 0666);
21     if (fout == -1) {
22         Fail("open out.txt");
23     }
24
25     // Replace stdin (fd 0) with our file fin
26     if (dup2(fin, 0) == -1) {
```

```
27     Fail("dup2 fin -> 0");
28 }
29
30 // Replace stdout (fd 1) with our file fout
31 if (dup2(fout, 1) == -1) {
32     Fail("dup2 fout -> 1");
33 }
34
35 // The original fds fin and fout can be closed,
36 // as their copies now exist as fd 0 and 1.
37 close(fin);
38 close(fout);
39 }
40
41 int main() {
42     Redirect();
43
44     int a;
45     std::cin >> a; // Now reads from input.txt
46     std::cout << a + 1 << std::endl; // Now writes to out.txt
47
48     return 0;
49 }
```

Листинг 2.3 – Функция перенаправления ввода-вывода

Если в файле `input.txt` будет число 123, то после выполнения программы в файле `out.txt` появится 124. Программа `main` ничего не знает о подмене; для неё `std::cin` и `std::cout` продолжают работать со стандартными [Файловый дескриптор](#) 0 и 1, но ядро теперь направляет эти операции в файлы.

Примечание

Проблемы буферизации. Стандартные потоки C++ (и C) буферизуют вывод для повышения производительности. Данные не отправляются ядру немедленно, а накапливаются во внутреннем буфере. Сброс буфера (`flush`) происходит:

- При его заполнении.
- При выводе специального символа, например, при использовании `std::endl`.
- При чтении из `std::cin` (обычно `stdout` сбрасывается).
- При завершении программы.

Интересно, что `libc` может менять свою стратегию буферизации. При выводе в терминал буфер часто сбрасывается при каждом символе новой строки (`'n'`). При выводе в файл (который не является интерактивным устройством) буферизация становится полной, и сброс происходит только при заполнении буфера или явном вызове `flush`. Это может приводить к неожиданному поведению, когда вывод, видимый в терминале, не сразу появляется в файле при перенаправлении.

2.5 Работа с директориями

Для просмотра содержимого директории используются функции из стандартной библиотеки C, которые являются обёрткой над соответствующими [Системный вызов](#).

```
1 #include <dirent.h>
2
3 DIR *opendir(const char *name);
4 struct dirent *readdir(DIR *dirp);
5 int closedir(DIR *dirp);
```

Листинг 2.4 – Интерфейс для чтения директорий

- `opendir` открывает директорию и возвращает указатель на структуру `DIR`, которая используется для дальнейших операций.
- `readdir` при каждом вызове возвращает указатель на структуру `dirent`, описывающую следующий элемент в директории. Когда элементы заканчиваются или происходит ошибка, возвращается `NULL`.
- `closedir` закрывает директорию.

Структура `dirent` содержит как минимум два поля: `d_name` (имя файла) и `d_type` (тип файла, например, `DT_REG` для файла, `DT_DIR` для директории).

```
1 #include <dirent.h>
2 #include <stdio.h>
3 #include <errno.h>
4
5 int main() {
6     DIR* dir = opendir(".");
7     if (!dir) {
8         perror("opendir failed");
9         return 1;
10    }
11
12    errno = 0; // To distinguish end-of-stream from error
13    struct dirent* entry;
14    while ((entry = readdir(dir)) != NULL) {
15        printf("%s\n", entry->d_name);
16    }
17
18    if (errno != 0) {
19        perror("readdir failed");
20    }
21
22    closedir(dir);
23    return 0;
24 }
```

Листинг 2.5 – Пример простой реализации `ls`

Примечание

В каждой директории в Linux есть два специальных вхождения:

- `..`: Ссылка на саму директорию.
- `...`: Ссылка на родительскую директорию.

Они также будут перечислены при вызове `readdir`.

Глава 3

3 Лекция

3.1 Дополнительные инструменты для работы с файловой системой

На прошлом занятии мы рассмотрели основы работы с файловой системой. Сегодня мы завершим эту тему, изучив несколько оставшихся, но важных инструментов, которые могут пригодиться в практических задачах.

3.1.1 Новые флаги для системного вызова `open`

Системный вызов `open` имеет несколько полезных флагов, которые мы не обсуждали ранее.

- `O_TRUNC`: Этот флаг позволяет при открытии файла немедленно обрезать его размер до нуля. Это удобная альтернатива последовательному вызову `open` и `ftruncate`, если содержимое файла нужно полностью перезаписать.
- `O_PATH`: Позволяет получить файловый дескриптор, который ссылается не на сам файл, а на его путь в файловой системе. Такой дескриптор имеет ограниченное применение (например, из него нельзя читать или в него писать), но он полезен для передачи в другие системные вызовы, такие как `fstat`, для получения информации об объекте файловой системы (включая директории), не открывая его для операций ввода-вывода.
- `O_NOFOLLOW`: Если путь, передаваемый в `open`, является символической ссылкой, то с этим флагом вызов не будет переходить по ней, а вернёт ошибку. Это важно для безопасности, чтобы избежать работы с непредусмотренным файлом.

3.1.2 Получение метаданных о файлах: семейство `stat`

Для получения подробной информации о файле или директории используется семейство системных вызовов `stat`. Они заполняют структуру `struct stat`, содержащую метаданные об объекте.

```
1 #include <sys/stat.h>
2
3 int stat(const char* path, struct stat* statbuf);
4 int lstat(const char* path, struct stat* statbuf);
5 int fstat(int fd, struct stat* statbuf);
```

Листинг 3.1 – Function signatures of the stat family

Ключевые различия между вызовами:

- `stat`: Принимает путь к файлу. Если путь указывает на символическую ссылку, `stat` переходит по ней и возвращает информацию о файле, на который она указывает.
- `lstat`: Аналогичен `stat`, но **не** переходит по символическим ссылкам. Вместо этого он возвращает информацию о самой ссылке.
- `fstat`: Принимает файловый дескриптор, полученный ранее через `open`.

Структура `struct stat` содержит множество полезных полей:

- `st_mode`: Тип файла (обычный файл, директория, символическая ссылка и т.д.) и права доступа к нему (чтение, запись, исполнение для владельца, группы и остальных).
- `st_uid` и `st_gid`: ID пользователя и группы-владельца файла.
- `st_size`: Размер файла в байтах.
- `st_blocks`: Количество дисковых блоков, занимаемых файлом.

- `st_atim`, `st_mtim`, `st_ctim`: Временные метки последнего доступа, последней модификации содержимого и последней модификации метаданных соответственно.

```

1  #include <fcntl.h>
2  #include <sys/stat.h>
3  #include <unistd.h>
4
5  // ...
6
7  int fd = open(path, O_RDONLY | O_PATH | O_NOFOLLOW);
8  if (fd == -1) { /* handle error */ }
9
10 struct stat stats;
11 if (fstat(fd, &stats) == -1) { /* handle error */ }
12
13 close(fd);
14 if (S_ISDIR(stats.st_mode)) {
15     // Directory
16 } else if (S_ISLNK(stats.st_mode)) {
17     // Symbol link
18 } else if (S_ISREG(stats.st_mode)) {
19     // Just file
20 }

```

Листинг 3.2 – Example of using `fstat` to determine the object type

В этом примере используется флаг `O_PATH`, чтобы безопасно получить дескриптор для проверки типа объекта, не открывая его для полноценной работы.

3.2 Управление памятью: виртуальная адресация

3.2.1 Проблема модели линейной памяти

Мы привыкли думать о памяти как о большом непрерывном массиве байтов. Однако эта модель не соответствует действительности. Проведём простой эксперимент: создадим две переменные — одну в *куча* (через `new`), а другую на *стек* (локальная переменная) — и выведем их адреса.

```

1  #include <iostream>
2
3  int main() {
4      int* heap_var = new int(10);
5      int stack_var = 20;
6
7      std::cout << "Heap address: " << (void*)heap_var << std::endl;
8      std::cout << "Stack address: " << (void*)&stack_var << std::endl;
9
10     long long diff = (long long)&stack_var - (long long)heap_var;
11     std::cout << "Difference (bytes): " << diff << std::endl;
12     // On a 64-bit system, the difference can be tens of terabytes
13
14     delete heap_var;
15     return 0;
16 }

```

Листинг 3.3 – Comparing stack and heap addresses

Разница между этими адресами может составлять десятки терабайт, что очевидно превышает объём физической оперативной памяти любого современного компьютера. Это наблюдение доказывает, что адреса, с которыми мы работаем в программе, не являются прямыми физическими адресами.

3.2.2 Виртуальная и физическая память

Для решения проблемы изоляции и безопасности процессов операционные системы вводят абстракцию — **виртуальная память**.

Определение: Виртуальная и физическая память

- **физическая память** — это реальные микросхемы оперативной памяти (RAM) в компьютере. Её адреса последовательны и ограничены её физическим объёмом.
- **виртуальная память** — это логическое адресное пространство, которое ОС предоставляет каждому процессу. Каждый процесс «видит» свой собственный, изолированный массив памяти, начинающийся с нуля. Адреса в этом пространстве называются **виртуальными**.

Процессор с помощью специального модуля (MMU — Memory Management Unit) и при содействии операционной системы преобразует виртуальные адреса в физические при каждом обращении к памяти. Это преобразование прозрачно для программиста.

3.2.3 Страничная организация памяти

Преобразование адресов происходит не для каждого байта в отдельности, а для блоков памяти фиксированного размера, называемых **страница памяти**.

Примечание

На большинстве современных систем (x86-64) размер страницы составляет 4 килобайта (4096 bytes, или 0x1000 в шестнадцатеричной системе). Узнать точный размер страницы в системе можно с помощью вызова `sysconf(_SC_PAGESIZE)`.

Операционная система поддерживает для каждого процесса таблицу страниц, которая устанавливает соответствие между страницами виртуальной и физической памяти.

При обращении к адресу, например, 0x2345:

1. Процессор разделяет его на номер страницы и смещение. Для страниц размером 0x1000 адрес 0x2345 — это смещение 0x345 внутри страницы 2.
2. С помощью таблицы страниц находится физический фрейм, соответствующий виртуальной странице 2 (на рис. 3.1 это фрейм А).
3. Процессор обращается к физической памяти по адресу, равному начальному адресу фрейма А плюс смещение 0x345.

Итоги раздела

Виртуальная память обеспечивает изоляцию процессов, позволяет программам работать с большим адресным пространством, чем доступно физической памяти,

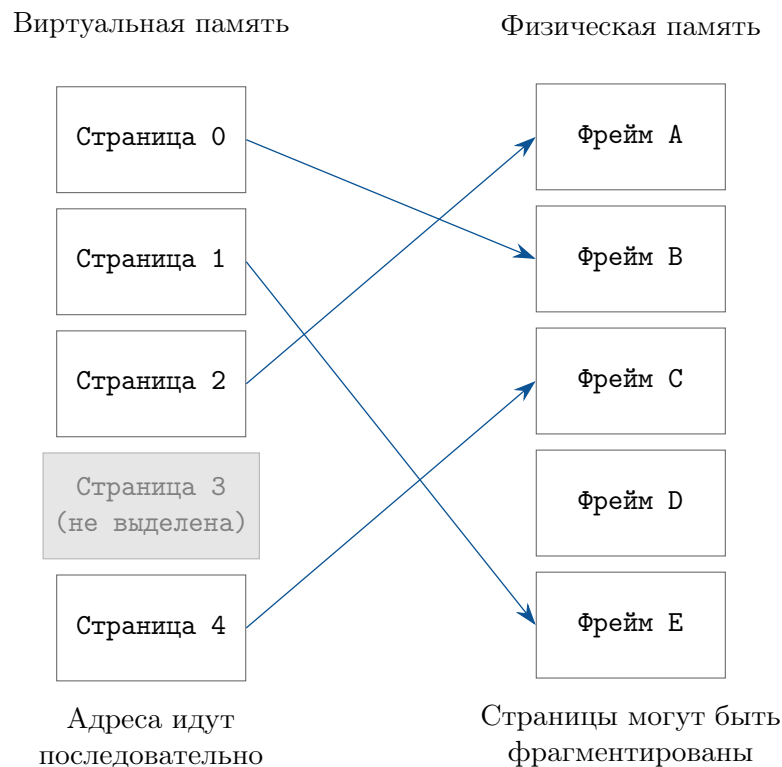


Рис. 3.1 – Схема отображения виртуальных страниц на физические фреймы памяти. Две соседние виртуальные страницы не обязательно отображаются в соседние физические.

и упрощает управление памятью для ОС. Это достигается за счёт постраничного отображения виртуальных адресов на физические.

3.3 Системные вызовы для управления памятью: `mmap`

Для управления виртуальным адресным пространством процесса в POSIX-системах используется системный вызов `mmap` и его пара `munmap`.

```

1 #include <sys/mman.h>
2
3 void *mmap(void *addr, size_t length, int prot, int flags,
4           int fd, off_t offset);
5 int munmap(void *addr, size_t length);

```

Листинг 3.4 – Signatures of `mmap` and `munmap`

`mmap` — это мощный, но сложный инструмент, который выполняет две основные функции:

1. **Анонимное отображение:** выделение новых страниц оперативной памяти для процесса.
2. **Файловое отображение:** отображение содержимого файла (или его части) в виртуальное адресное пространство процесса.

3.3.1 Аргументы и флаги `mmap`

Рассмотрим ключевые параметры `mmap`:

- **addr**: Желаемый стартовый адрес для отображения. Обычно передаётся `nullptr`, чтобы ОС сама выбрала подходящий адрес.
- **length**: Размер отображаемой области в байтах.
- **prot** (protection): Права доступа к памяти.
 - `PROT_READ`: память можно читать.
 - `PROT_WRITE`: в память можно писать.
 - `PROT_EXEC`: содержимое памяти можно исполнять как код.
 - `PROT_NONE`: к памяти нет доступа.
- **flags**: Определяют тип и поведение отображения.
 - `MAP_SHARED` или `MAP_PRIVATE`: Один из этих флагов обязателен. `MAP_SHARED` означает, что изменения, сделанные в памяти, будут видны другим процессам, отображающим тот же объект, и (в случае файла) будут записаны обратно в файл. `MAP_PRIVATE` создаёт сору-он-вайт отображение: изменения видны только текущему процессу и не затрагивают исходный файл.
 - `MAP_ANONYMOUS`: Создаёт анонимное отображение. Память инициализируется нулями и не связана ни с каким файлом. При использовании этого флага аргумент `fd` должен быть `-1`.
 - `MAP_FIXED`: Требует от ОС использовать точно адрес, указанный в `addr`. Это опасный флаг, так как он может без предупреждения перезаписать существующие отображения.
- **fd, offset**: Файловый дескриптор и смещение от начала файла для файловых отображений.

В случае успеха `mmap` возвращает указатель на начало выделенной области. В случае ошибки — `MAP_FAILED`.

3.3.2 Примеры использования `mmap`

Анонимное отображение

Это основной способ, которым аллокаторы (`malloc`, `new`) запрашивают большие блоки памяти у операционной системы.

```
1 #include <sys/mman.h>
2 #include <unistd.h> // For sysconf
3
4 // ...
5
6 // Request one page of memory
7 size_t page_size = sysconf(_SC_PAGESIZE);
8 void* raw_mem = mmap(nullptr, page_size,
9                       PROT_READ | PROT_WRITE,
10                      MAP_PRIVATE | MAP_ANONYMOUS,
11                      -1, 0);
12 if (raw_mem == MAP_FAILED) {
13     // Error handling
14 }
15
```

```
16 char* data = static_cast<char*>(raw_mem);
17 // Now 'data' can be used as a regular array
18 data[0] = 'H';
19 data[1] = 'i';
20
21 // Free the memory
22 munmap(raw_mem, page_size);
```

Листинг 3.5 – Allocating one page of memory using mmap

Отображение файла в память

Отображение файла позволяет работать с его содержимым как с обычным массивом в памяти, что может быть эффективнее, чем многократные вызовы `read` и `write`, особенно при произвольном доступе.

```
1 #include <sys/mman.h>
2 #include <fcntl.h>
3 #include <unistd.h>
4
5 const size_t FILE_SIZE = 128;
6 int fd = open("storage.bin", O_RDWR | O_CREAT, 0644);
7 ftruncate(fd, FILE_SIZE);
8 // Set the file size
9
10 void* raw_mem = mmap(nullptr, FILE_SIZE,
11                      PROT_READ | PROT_WRITE,
12                      MAP_SHARED, // Changes will be written to the file
13                      fd, 0);
14 close(fd); // The file descriptor can be closed after mmap
15
16 if (raw_mem == MAP_FAILED) { /* ... */ }
17
18 char* data = static_cast<char*>(raw_mem);
19 for (size_t i = 0; i < FILE_SIZE; ++i) {
20     data[i] = static_cast<char>(i);
21 }
22
23 // The operating system will write the changes to the disk
24 // (not necessarily immediately)
25
26 munmap(raw_mem, FILE_SIZE);
```

Листинг 3.6 – Working with a file via mmap

3.3.3 Освобождение памяти: `munmap`

Вызов `munmap` удаляет отображение для указанного диапазона виртуальных адресов. Крайне важно освобождать память, выделенную через `mmap`, чтобы избежать утечек ресурсов. Аналогично паре `new/delete`, каждому успешному вызову `mmap` должен соответствовать вызов `munmap`.

3.4 Аргументы командной строки и переменные окружения

Кроме ввода-вывода, программа может получать информацию извне при запуске. Рассмотрим два основных механизма: аргументы командной строки и переменные окружения.

3.4.1 Аргументы командной строки

При запуске программы из терминала можно передать ей параметры. Они доступны в функции `main` через её аргументы.

```
1 int main(int argc, char* argv[]) {  
2     // ...  
3 }
```

Листинг 3.7 – The main function interface

- `argc` (argument count): количество переданных аргументов.
- `argv` (argument vector): массив указателей на C-строки.

Важно помнить, что `argv[0]` — это всегда имя самой запущенной программы. Реальные аргументы начинаются с `argv[1]`. Например, для команды `./myprog hello world` будет:

- `argc = 3`
- `argv[0] = "./myprog"`
- `argv[1] = "hello"`
- `argv[2] = "world"`

3.4.2 Переменные окружения

Переменные окружения — это набор пар "ключ-значение" которые наследуются дочерними процессами от родительских. Они используются для передачи контекста и настроек программам (например, `PATH` для поиска исполняемых файлов, `HOME` для пути к домашней директории). В Linux переменные окружения физически располагаются в памяти процесса сразу после массива `argv`, отделённые от него указателем `nullptr`. Для безопасного доступа к ним из C++ используется функция `getenv`.

```
1 #include <iostream>  
2 #include <cstdlib> // For getenv  
3  
4 int main() {  
5     const char* user = std::getenv("USER");  
6     if (user != nullptr) {  
7         std::cout << "Hello, " << user << "!"  
8         << std::endl;  
9     } else {  
10        std::cout << "USER environment variable is not set."  
11        << std::endl;  
12    }  
13    return 0;  
14 }
```

Листинг 3.8 – Reading an environment variable

Примечание

Переменные окружения часто используются для передачи конфиденциальной информации (ключей API, паролей), так как они, в отличие от аргументов командной строки, не видны другим пользователям системы через команды типа `ps`.

Итоги раздела

- Аргументы командной строки (`argc`, `argv`) позволяют передавать простые параметры при запуске.
- Переменные окружения — это наследуемые пары "ключ-значение" для передачи настроек и контекста.
- Для доступа к переменным окружения следует использовать `getenv`, что является более безопасным и портируемым способом.

Глава 4

4 Лекция

4.1 Углублённая работа с памятью

На прошлой лекции мы познакомились с концепцией **виртуальная память** и системным вызовом `mmap`, который управляет отображением виртуальных адресов на **физическая память**. Однако модель, в которой `mmap` немедленно выделяет реальные физические страницы, является упрощением. На практике современные ОС используют более сложный и эффективный механизм.

4.1.1 Механизм Page Fault и ленивое выделение памяти

При попытке программы обратиться по виртуальному адресу, который не сопоставлен ни одной физической странице, процессор генерирует специальное прерывание.

Определение: Страничная ошибка (Page Fault)

страничная ошибка — это прерывание, которое генерируется аппаратно (процессором) при попытке доступа к странице **виртуальная память**, не имеющей корректного отображения в **физическая память**. При возникновении **страничная ошибка** исполнение текущего кода программы приостанавливается, и управление передаётся обработчику в ОС.

ОС анализирует причину **страничная ошибка**. Если обращение было к некорректному адресу (например, разыменован нулевой указатель), ОС принудительно завершает программу, как правило, с ошибкой **Segmentation Fault**.

Однако этот же механизм используется для реализации **ленивого выделения памяти** (**постраничная подкачка по требованию**). Когда программа вызывает `mmap`, ОС на самом деле не выделяет физические страницы. Она лишь запоминает, что данный диапазон виртуальных адресов теперь является валидным для процесса. Реальное выделение физической страницы происходит только при **первом обращении** к ней. Это обращение вызывает **страничная ошибка**, который ОС обрабатывает:

1. Находит свободную физическую страницу.
2. Устанавливает отображение между виртуальной страницей, вызвавшей прерывание, и новой физической страницей.
3. Возобновляет исполнение программы с прерванной инструкции.

Для программы этот процесс прозрачен, за исключением небольшой задержки.

Примечание

Плюсы и минусы ленивого выделения:

- **Плюс:** Эффективное использование ресурсов. Программы часто запрашивают больше памяти, чем реально используют. Ленивый подход позволяет системе выделять только ту физическую память, которая действительно нужна, и поддерживать так называемый *memory overcommitment* (когда суммарный объем запрошенной памяти превышает имеющуюся физическую).
- **Минус:** Усложнение обработки ошибок нехватки памяти. Вместо проверки кода возврата `mmap`, программа может быть внезапно "убита" ОС в произвольный момент при обращении к памяти, если свободные физические страницы закончились.
- **Минус:** Непредсказуемые задержки. Обращение к "новой" странице памяти вызывает **страничная ошибка**, что приводит к задержке, так как управление передается

ОС. Это может быть критично для приложений реального времени.

4.1.2 Структура таблиц страниц (Page Tables)

Для трансляции виртуальных адресов в физические ОС и процессор используют **таблица страниц**. Хранить простое линейное отображение для всего 64-битного адресного пространства (даже с учётом реальных ограничений современных процессоров в 256 ТБ) неэффективно.

В архитектуре x86-64 используется **четырёхуровневая древовидная структура** таблиц страниц. Виртуальный адрес делится на несколько частей:

- **Смещение (offset):** Младшие 12 бит, указывающие на байт внутри страницы ($2^{12} = 4096$ байт).
- **Индексы в таблицах:** Четыре группы по 9 бит каждая, которые используются для последовательного обхода четырёхуровневого дерева таблиц (L3, L2, L1, L0).

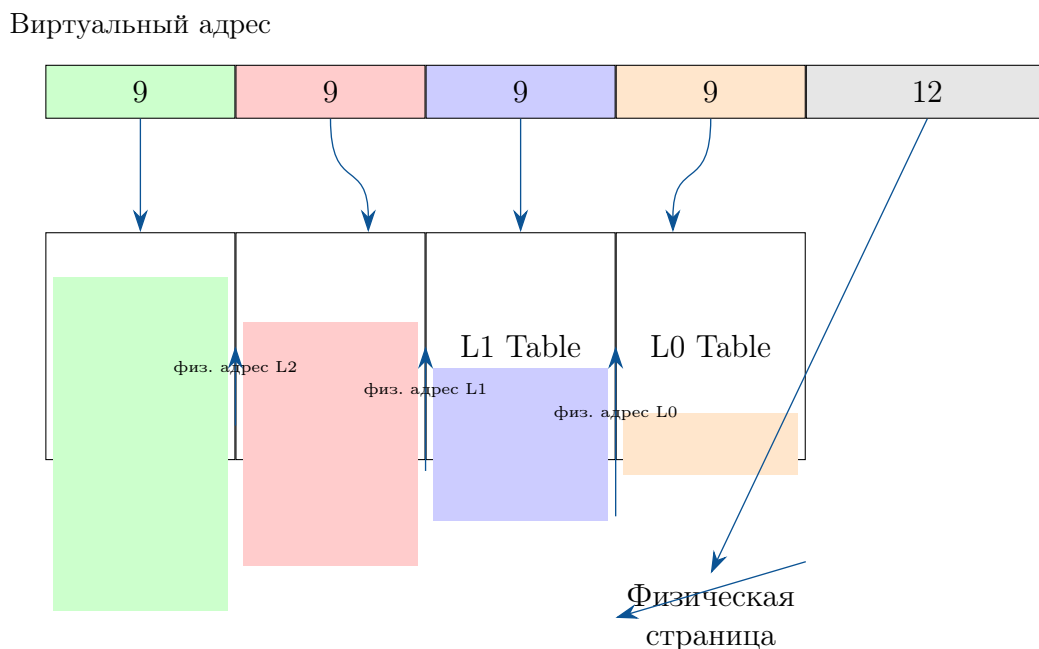


Рис. 4.1 – Трансляция виртуального адреса в физический через четырёхуровневые таблицы страниц.

Процессор аппаратно выполняет обход этой структуры при каждом доступе к памяти: использует 9 бит адреса как индекс в таблице L3, находит там физический адрес таблицы L2, затем следующие 9 бит — как индекс в L2, и так далее, пока не дойдет до таблицы L0, где хранится адрес искомой физической страницы.

4.1.3 Дополнительные системные вызовы для работы с памятью

`mprotect(addr, size, prot)` изменяет права доступа (чтение, запись, исполнение) для уже выделенного диапазона виртуальной памяти `[addr, addr+size)`.

`mremap(old_addr, old_size, new_size, flags, ...)` позволяет изменять размер существующего отображения, а также перемещать его на новое место в виртуальном адресном пространстве.

`mlock(addr, size)` "закрепляет" указанный диапазон страниц в физической памяти,

запрещая ОС выгружать их в **своп (swap)**. Это важно для приложений, работающих с чувствительными данными (пароли, ключи шифрования) или требующих предсказуемых задержек. `munlock` отменяет это действие.

Особый режим `mremap` позволяет создать "копию" участка памяти, где два разных диапазона виртуальных адресов указывают на **одни и те же физические страницы**. Любая запись в один диапазон немедленно видна в другом.

```
1 // Allocate original mapping
2 void* from = mmap(nullptr, PAGE_SIZE, PROT_READ | PROT_WRITE,
3                   MAP_ANONYMOUS | MAP_SHARED, -1, 0);
4
5 // Reserve space for the "copy"
6 void* to_placeholder = mmap(nullptr, PAGE_SIZE, PROT_NONE,
7                              MAP_ANONYMOUS | MAP_PRIVATE, -1, 0);
8
9 // Create the shared mapping (remap `from` onto `to_placeholder`)
10 // MREMAP_FIXED tells mremap to use the address we provide.
11 // old_size = 0 is a special value for this copy operation.
12 void* to = mremap(from, 0 /* old_size */, PAGE_SIZE,
13                  MREMAP_MAYMOVE | MREMAP_FIXED, to_placeholder);
14
15 // Now, `from` and `to` point to the same physical page.
16 volatile int* p_from = static_cast<volatile int*>(from);
17 volatile int* p_to = static_cast<volatile int*>(to);
18
19 *p_from = 123;
20 // Reading from p_to will now yield 123.
21 printf("Value at 'to': %d\n", *p_to); // Prints 123
```

Листинг 4.1 – Пример использования `mremap` для создания разделяемого отображения.

Примечание

В листинг 4.1 используется ключевое слово `volatile`. Оно сообщает компилятору, что значение в памяти, на которую указывает указатель, может измениться в любой момент без его ведома (например, через другой указатель, как в нашем случае). Это запрещает компилятору кэшировать значение переменной в регистре и заставляет его каждый раз честно читать значение из памяти, предотвращая неверные оптимизации.

4.1.4 Swap (своп) и его проблемы

Когда физическая память заканчивается, ОС может использовать **своп (swap)**: выгрузить содержимое некоторых "неактивных" физических страниц на жесткий диск, чтобы освободить место для более актуальных данных. Когда программа обратится к такой выгруженной странице, произойдет **страничная ошибка**, и ОС загрузит её обратно с диска.

Проблемы свопинга:

- **Производительность:** Диск значительно медленнее оперативной памяти, что приводит к большим задержкам.
- **Безопасность:** Секретные данные (ключи, пароли) могут оказаться на диске в незашифрованном виде и остаться там даже после выключения питания, создавая уязви-

мость.

Итоги раздела

- Обращение к неотмеченной в [таблица страниц](#) странице вызывает [страничная ошибка](#).
- ОС использует [страничная ошибка](#) для реализации ленивого выделения памяти, что экономит физическую память.
- Трансляция адресов в x86-64 реализована через многоуровневые таблицы страниц.
- Системные вызовы `mprotect`, `mremap`, `mlock` предоставляют тонкий контроль над отображениями памяти.
- [своп \(swap\)](#) помогает при нехватке памяти, но ценой производительности и потенциальных рисков безопасности.

4.2 Управление процессами

До сих пор мы рассматривали работу в рамках одного процесса. Теперь изучим, как создавать новые процессы и управлять ими.

Определение: Процесс

Процесс — это экземпляр запущенной программы. Каждый процесс является изолированной сущностью и обладает собственными ресурсами:

- Уникальным **Process ID** (идентификатор процесса) (PID).
- Отдельным виртуальным адресным пространством.
- Собственной таблицей файловых дескрипторов.

Процессы могут выполняться параллельно на многоядерных системах.

4.2.1 Подмена процесса: семейство `exec`

Системные вызовы семейства `exec` (`execlp`, `execvpe` и др.) **не создают** новый процесс. Они полностью **заменяют** текущий процесс новым, загружая и запуская указанный исполняемый файл.

```
1 #include <unistd.h>
2 #include <stdio>
3
4 int main() {
5     printf("Before exec...\n");
6
7     // Replace the current process with "ls -l"
8     // The first argument is the command,
9     // subsequent args are for its argv.
10    // The list must be terminated by a NULL pointer.
11    execlp("ls", "ls", "-l", nullptr);
12
13    // This line will never be reached if execlp succeeds.
14    perror("execlp failed");
15    return 1;
16 }
```

Листинг 4.2 – Запуск утилиты `ls` с помощью `execlp`.

При успешном вызове `execlp` код после него никогда не выполняется. Новый процесс (в данном случае, `ls`) наследует некоторые атрибуты старого, например, таблицу файловых дескрипторов, но получает новое адресное пространство.

Примечание

При запуске сторонних программ важно избегать утечки файловых дескрипторов. Если библиотека внутри вашего кода открыла файл, он останется открытым и в запущенном через `exec` процессе. Стандартное решение — открывать все файловые дескрипторы с флагом `O_CLOEXEC`, который предписывает ядру автоматически закрыть этот дескриптор при вызове `exec`.

4.2.2 Создание процесса: `fork`

Для создания нового процесса используется системный вызов `fork`.

Определение: Системный вызов fork

`fork()` создаёт точную копию текущего процесса. Уникальность `fork` в том, что он вызывается один раз, а возвращается дважды:

- В родительском процессе `fork()` возвращает `PID` нового (дочернего) процесса.
- В дочернем процессе `fork()` возвращает `0`.
- В случае ошибки возвращается `-1`.

Дочерний процесс является почти полной копией родителя: он получает копию адресного пространства, стека вызовов и таблицы файловых дескрипторов. Исполнение в обоих процессах продолжается с точки сразу после вызова `fork`.

```
1 #include <unistd.h>
2 #include <sys/wait.h>
3 #include <stdio.h>
4
5 int main() {
6     pid_t child_pid = fork();
7
8     if (child_pid == -1) {
9         perror("fork failed");
10        return 1;
11    } else if (child_pid == 0) {
12        // We are in the child process
13        printf("I am the child! My PID is %d\n", getpid());
14    } else {
15        // We are in the parent process
16        printf("I am the parent! My child's PID is %d\n", child_pid);
17        wait(nullptr); // Wait for the child to finish
18        printf("Parent knows child has finished.\n");
19    }
20
21    return 0;
22 }
```

Листинг 4.3 – Базовое использование `fork`.

4.2.3 Жизненный цикл процесса

Завершение процесса: `exit` vs `_Exit`

- `std::exit(code)` — функция стандартной библиотеки. Она не только завершает процесс с кодом `code`, но и выполняет ряд "очищающих" действий: сбрасывает буферы потоков ввода-вывода (например, `cout`), вызывает обработчики, зарегистрированные через `atexit`, и т.д..
- `std::_Exit(code)` (или системный вызов `_exit`) — немедленно завершает процесс без какой-либо очистки. В дочерних процессах после `fork` предпочтительнее использовать именно `_Exit`, чтобы избежать нежелательных побочных эффектов, например, двойного сброса буферов, которые были скопированы от родителя.

Ожидание дочерних процессов: `wait` и `waitpid`

Родительский процесс обязан "собирать" информацию о завершении своих дочерних процессов с помощью `wait()` или `waitpid()`. Эти вызовы блокируют родителя до тех пор, пока один из его детей не завершится, и позволяют получить его код завершения.

Определение: Процесс-зомби

процесс-зомби — это процесс, который уже завершил своё выполнение, но запись о нём (PID, код завершения) всё ещё остаётся в таблице процессов ядра. Он находится в этом состоянии до тех пор, пока родитель не "прочитает" его статус с помощью `wait`. Если родитель не делает `wait`, зомби накапливаются и "утекают" системные ресурсы (в частности, PID).

Определение: Процесс-сирота

процесс-сирота — это процесс, родитель которого завершился раньше него. Такие процессы не остаются "бесхозными" — их "усыновляет" специальный системный процесс `init` (с PID 1), который периодически вызывает `wait` и очищает зомби.

4.2.4 Паттерн `fork-exec`

Комбинация `fork` и `exec` — это стандартный способ в Unix-системах запустить новую программу, не прекращая работу текущей.

1. Родительский процесс вызывает `fork()`, создавая свою копию.
2. В дочернем процессе (где `fork()` вернул 0) выполняются необходимые настройки (например, перенаправление ввода-вывода с помощью `dup2`).
3. Дочерний процесс вызывает один из вызовов семейства `exec`, заменяя себя новой программой.
4. Родительский процесс (где `fork()` вернул `PID > 0`) может продолжить свою работу или дождаться завершения дочернего с помощью `waitpid()`.

Итоги раздела

- Процесс — это изолированный экземпляр запущенной программы.
- `exec` заменяет текущий процесс, `fork` создаёт его копию.
- Паттерн `fork-exec` является основой для запуска программ в Unix-подобных системах.
- Родитель **обязан** дожидаться завершения дочерних процессов с помощью `wait` или `waitpid`, чтобы избежать появления **процесс-зомби**.

4.3 Межпроцессное взаимодействие: Pipelines

Одним из самых мощных механизмов в Unix является **канал (pipe)**, который позволяет связать стандартный вывод одного процесса со стандартным вводом другого. Рассмотрим, как реализовать аналог команды `ps aux | grep zsh` программно.

Для этого нам понадобится системный вызов `pipe()`, который создаёт однонаправленный канал данных и возвращает два файловых дескриптора: `pipefd[0]` для чтения и `pipefd[1]` для записи.

Алгоритм реализации пайплайна `cmd1 | cmd2`:

1. Создать **канал (pipe)** с помощью `pipe(pipefd)`.
2. Вызвать `fork()` для создания первого дочернего процесса (`child1` для `cmd1`).
3. В `child1`:
 - Закрыть ненужный конец канала: `close(pipefd[0])`.
 - Перенаправить стандартный вывод на пишущий конец канала: `dup2(pipefd[1], STDOUT_FILENO)`.
 - Закрыть оригинальный дескриптор: `close(pipefd[1])`.
 - Вызвать `exec` для запуска `cmd1`.
4. Вызвать `fork()` для создания второго дочернего процесса (`child2` для `cmd2`).
5. В `child2`:
 - Закрыть ненужный конец канала: `close(pipefd[1])`.
 - Перенаправить стандартный ввод на читающий конец канала: `dup2(pipefd[0], STDIN_FILENO)`.
 - Закрыть оригинальный дескриптор: `close(pipefd[0])`.
 - Вызвать `exec` для запуска `cmd2`.
6. В родительском процессе:
 - **Критически важно:** закрыть **оба** конца канала: `close(pipefd[0])` и `close(pipefd[1])`. Если этого не сделать, читающий процесс никогда не получит EOF и зависнет.
 - Дождаться завершения обоих дочерних процессов с помощью `waitpid()`.

```
1 #include <unistd.h>
2 #include <sys/wait.h>
3 #include <stdio.h>
4 #include <stdlib.h>
5
6 int main() {
7     int pipefd[2];
8     if (pipe(pipefd) == -1) {
9         perror("pipe failed");
10        return 1;
11    }
12
13    pid_t child1 = fork();
14    if (child1 == 0) { // First child: ps aux
15        close(pipefd[0]); // Close read end
```

```
16     dup2(pipefd[1], STDOUT_FILENO);
17     close(pipefd[1]); // Close original write end
18     execlp("ps", "ps", "aux", nullptr);
19     _exit(1); // exit if exec fails
20 }
21
22 pid_t child2 = fork();
23 if (child2 == 0) { // Second child: grep zsh
24     close(pipefd[1]); // Close write end
25     dup2(pipefd[0], STDIN_FILENO);
26     close(pipefd[0]); // Close original read end
27     execlp("grep", "grep", "zsh", nullptr);
28     _exit(1); // exit if exec fails
29 }
30
31 // Parent process
32 close(pipefd[0]); // ESSENTIAL: close both pipe ends in parent
33 close(pipefd[1]);
34
35 waitpid(child1, nullptr, 0);
36 waitpid(child2, nullptr, 0);
37
38 return 0;
39 }
```

Листинг 4.4 – Программная реализация пайплайна `ps aux | grep zsh`.

Примечание

Что происходит, если читатель завершается раньше писателя? (например, в `ps aux | head -n 5`) Когда все читающие концы канала закрываются, а писатель пытается в него что-то записать, ОС посылает писателю сигнал **SIGPIPE**. По умолчанию, действие для этого сигнала — аварийное завершение процесса. Это элегантно решает проблему "бесконечной" работы процессов в начале пайплайна, если их вывод больше никому не нужен.

Итоги раздела

- **канал (pipe)** создаёт однонаправленный канал для данных между процессами.
- Комбинация `pipe`, `fork`, `dup2` и `exec` позволяет строить сложные конвейеры обработки данных.
- В родительском процессе необходимо закрывать оба конца канала, чтобы избежать взаимоблокировок.
- Сигнал **SIGPIPE** автоматически завершает процессы, которые пытаются писать в "сломанный" канал (без читателей).

Глава 5

5 Лекция

5.1 Управление процессами: группы и сигналы

Продолжая изучение процессов как единиц параллелизма с изолированными адресными пространствами, мы рассмотрим механизмы их организации и взаимодействия. Для эффективного управления множеством связанных процессов операционные системы, включая Linux, предоставляют инструменты для их группировки и асинхронного уведомления.

5.1.1 Группы процессов (Process Groups)

Для упрощения управления несколькими процессами одновременно они могут быть объединены в группы. Каждый процесс в системе принадлежит определённой группе.

Определение: Группа процессов

PGID (Process Group ID) — это числовой идентификатор, общий для нескольких процессов. Он позволяет применять операции, такие как отправка сигналов, ко всей группе сразу, а не к каждому процессу по отдельности. Каждый процесс также принадлежит **SID (Session ID)**, которая объединяет группы процессов.

Для работы с **PGID (Process Group ID)** существуют системные вызовы:

- **getpgid(pid_t pid)**: получает **PGID (Process Group ID)** процесса с указанным **pid**. Вызов **getpgid(0)** вернёт **PGID (Process Group ID)** текущего процесса.
- **setpgid(pid_t pid, pid_t pgid)**: устанавливает **PGID (Process Group ID)** для процесса. Чтобы создать новую группу, обычно процесс вызывает **setpgid(0, 0)**, что создаёт новую группу с **PGID (Process Group ID)**, равным **PID** этого процесса.

5.1.2 Сигналы (Signals)

Определение: Сигнал

сигнал — это простой механизм **IPC (Inter-Process Communication)**, представляющий собой асинхронное уведомление, которое может быть отправлено процессу операционной системой или другим процессом. В отличие от пайпов, для отправки сигнала не требуется наличие родственной связи между процессами.

Сигналы, как и **PID**, являются просто числами. При получении сигнала процесс может отреагировать одним из нескольких способов:

- **Term (Termination)**: Завершение процесса. Это действие по умолчанию для большинства сигналов.
- **Ign (Ignore)**: Игнорирование сигнала.
- **Core**: Завершение процесса с генерацией **core dump**. Это файл, содержащий полный снимок адресного пространства процесса в момент сбоя, что позволяет проводить посмертную отладку (post-mortem debugging) с помощью таких инструментов, как GDB.
- **Stop**: Приостановка выполнения процесса.
- **Cont (Continue)**: Возобновление выполнения приостановленного процесса.

Процесс может переопределить стандартную реакцию на большинство сигналов, установив собственный обработчик.

Таблица 5.1 – Некоторые распространённые сигналы и их действия по умолчанию

Сигнал	Номер	Действие	Комментарий
SIGINT	2	Term	Отправляется при нажатии Ctrl+C в терминале.
SIGQUIT	3	Core	Отправляется при нажатии Ctrl+\ .
SIGKILL	9	Term	Гарантированно завершает процесс. Этот сигнал нельзя перехватить или проигнорировать.
SIGSEGV	11	Core	Segmentation Fault. Отправляется при попытке доступа к неразрешённой области памяти.
SIGTSTP	20	Stop	Отправляется при нажатии Ctrl+Z в терминале, приостанавливая процесс.

5.1.3 Отправка сигналов и ожидание процессов

Для отправки сигналов используется системный вызов `kill`. Несмотря на название, он может отправлять любой сигнал, а не только те, что завершают процесс.

```

1 #include <signal.h>
2 int kill(pid_t pid, int sig);

```

Листинг 5.1 – Сигнатура системного вызова `kill`

Аргумент `pid` интерпретируется следующим образом:

- `pid > 0`: Сигнал `sig` отправляется процессу с PID, равным `pid`.
- `pid < -1`: Сигнал отправляется всем процессам в группе с **PGID (Process Group ID)**, равным `-pid`.
- `pid == 0`: Сигнал отправляется всем процессам в группе текущего процесса.
- `pid == -1`: Сигнал отправляется всем процессам, которым текущий пользователь имеет право отправлять сигналы (за исключением некоторых системных процессов).

Механизмы ожидания дочерних процессов, такие как `wait` и `waitpid`, позволяют не только дождаться завершения, но и получить информацию о причине.

- `WIFEXITED(status)`: Возвращает `true`, если процесс завершился штатно через вызов `exit()`. Код возврата можно получить с помощью `WEXITSTATUS(status)`.
- `WIFSIGNALED(status)`: Возвращает `true`, если процесс был завершён сигналом. Номер сигнала можно получить через `WTERMSIG(status)`.

Вызов `waitpid` также интегрирован с группами процессов:

- `waitpid(-1, ...)`: Ждёт любого дочернего процесса (стандартное поведение).
- `waitpid(-pgid, ...)`: Ждёт завершения любого дочернего процесса из группы с **PGID (Process Group ID)**.

5.1.4 Управление памятью при `fork()`

Ранее мы говорили, что `fork()` создаёт полную копию адресного пространства родителя для дочернего процесса. Для современных процессов, занимающих гигабайты памяти, полное копирование было бы крайне неэффективным.

Определение: Copy-on-Write (COW)

Copy-on-Write (копирование при записи) — это оптимизация, применяемая при вызове `fork()`. Вместо реального копирования всех страниц памяти, операционная система создаёт для дочернего процесса новые таблицы страниц, которые указывают на те же физические страницы, что и у родителя. Все эти страницы помечаются как доступные только для чтения. При попытке записи в такую страницу (и родителем, и ребёнком) происходит аппаратное прерывание. ОС перехватывает его, создаёт реальную копию только этой конкретной страницы, и уже в неё производится запись. Это позволяет копировать только те данные, которые действительно изменяются, экономя время и память.

Общая память через `mmap`

Хотя **Copy-on-Write (копирование при записи)** обеспечивает изоляцию, иногда процессам нужна общая область памяти для эффективного взаимодействия. Этого можно достичь с помощью системного вызова `mmap` с флагом `MAP_SHARED`.

Примечание

Если участок памяти был создан с флагом `MAP_PRIVATE`, то при `fork()` к нему применяется **Copy-on-Write (копирование при записи)**. Если же был использован флаг `MAP_SHARED`, то этот участок памяти будет общим для родителя и ребёнка после `fork()`: изменения, сделанные одним процессом, будут видны другому.

Важно помнить, что при работе с общей памятью возникает проблема синхронизации. Нет гарантий относительно порядка выполнения инструкций в родителе и ребёнке. Чтобы корректно читать данные, записанные другим процессом, необходимо использовать механизмы синхронизации, например, пайпы или сигналы, чтобы уведомить читающий процесс о готовности данных.

Итоги раздела

- **Группы процессов (PGID (Process Group ID))** упрощают управление множеством процессов, позволяя применять операции ко всей группе сразу.
- **Сигналы** — это механизм асинхронных уведомлений для межпроцессного взаимодействия.
- Вызов `kill` позволяет отправлять сигналы процессам и группам процессов.
- `waitpid` может ожидать завершения процессов из определённой группы.
- **Copy-on-Write (Copy-on-Write (копирование при записи))** оптимизирует `fork()`, откладывая реальное копирование страниц памяти до первой операции записи.
- `mmap` с флагом `MAP_SHARED` создаёт общую область памяти для родителя и дочерних процессов, но требует явной синхронизации доступа.

5.2 Представление данных

Любые данные в компьютерных системах — будь то файлы, сетевые пакеты или потоки ввода-вывода — в конечном счёте представляются в виде последовательности байт. Способ преобразования структурированных данных в байты и обратно определяет формат данных.

5.2.1 Текстовые и бинарные форматы

Форматы данных условно делятся на две большие категории.

- **Текстовые форматы** (JSON, YAML, XML, TXT) оптимизированы для чтения и редактирования человеком. Они, как правило, избыточны и требуют больше ресурсов для парсинга (синтаксического анализа).
- **Бинарные форматы** (исполняемые файлы ELF, архивы ZIP, форматы сериализации BSON, Protocol Buffers) оптимизированы для машинной обработки, компактности или скорости. Они нечитаемы для человека без специальных инструментов, но обрабатываются программами гораздо эффективнее.

Некоторые бинарные форматы, например, исполняемые файлы ELF (Executable and Linkable Format), спроектированы так, что их можно загрузить в память простым отображением файла с помощью `mmap`, без дополнительной обработки.

Тонкости бинарных форматов: порядок байтов

При работе с бинарными данными, содержащими числа размером более одного байта, возникает проблема [порядок байтов \(endianness\)](#).

Определение: Порядок байтов (Endianness)

[порядок байтов \(endianness\)](#) определяет, как байты многобайтового числа располагаются в оперативной памяти.

- **Little-endian:** Младший байт числа хранится по младшему адресу памяти. Этот порядок используется в большинстве современных архитектур, включая x86-64.
- **Big-endian:** Старший байт числа хранится по младшему адресу. Этот порядок часто используется в сетевых протоколах (network byte order).

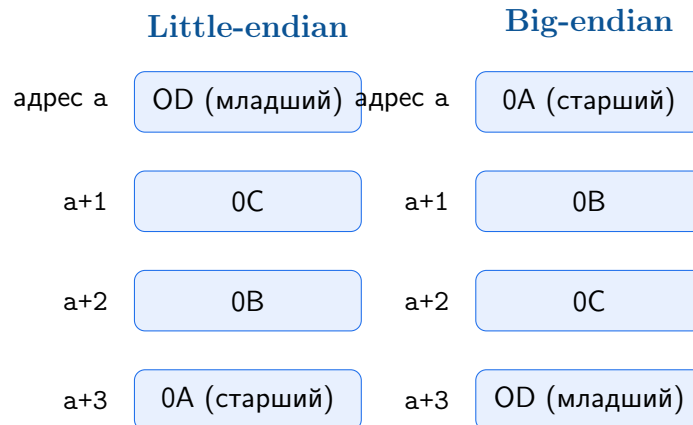
Проблема возникает при передаче бинарных данных между системами с разным порядком байтов. Если данные записываются и читаются на одной и той же машине, беспокоиться о порядке байтов не нужно.

5.2.2 Кодировки текста: от ASCII до Unicode

Представление текста — одна из фундаментальных задач. Исторически первой и наиболее влиятельной кодировкой стала [ASCII \(American Standard Code for Information Interchange\)](#).

Определение: ASCII

[ASCII \(American Standard Code for Information Interchange\)](#) — стандарт, определяющий соответствие между числами (кодами от 0 до 127) и символами. Он использует только 7 бит, восьмой бит всегда равен нулю. ASCII включает в себя символы английского алфавита, цифры, знаки препинания и управляющие символы.

32-битное число: 0x0A0B0C0D**Рис. 5.1** – Расположение байтов числа 0x0A0B0C0D в памяти при разном порядке байтов

Среди управляющих символов ASCII есть:

- 0x0A (`\n`) — перевод строки (Line Feed).
- 0x0D (`\r`) — возврат каретки (Carriage Return).
- 0x1B (ESC) — Escape, используется для начала управляющих последовательностей, например, для управления цветом текста в терминале.

Примечание

Исторически сложилось два основных способа кодирования конца строки в текстовых файлах:

- **Unix/Linux:** используется один символ `\n`.
- **DOS/Windows:** используется последовательность из двух символов `\r\n` (возврат каретки и перевод строки).

Это различие может вызывать проблемы, например, при запуске скриптов с Windows-окончаниями строк в Linux, так как интерпретатор может неверно обработать `\r` в конце строки shebang.

Очевидным недостатком [ASCII \(American Standard Code for Information Interchange\)](#) является отсутствие поддержки символов других языков. Это привело к появлению множества несовместимых 8-битных кодировок (семейства ISO-8859, CP1251, KOI8-R и др.), которые использовали старший бит для кодирования национальных алфавитов. Проблема усугублялась в языках с иероглифической письменностью, где требовались многобайтовые кодировки со сложной логикой переключения режимов (например, EUC-JP).

5.2.3 Unicode и его кодировки

Для решения проблемы хаоса кодировок был создан стандарт [Unicode](#).

Определение: Unicode

Unicode — это стандарт, который сопоставляет каждому символу из большинства мировых письменностей (включая мёртвые языки и эмодзи) уникальное целое число, называемое **Unicode**. Всего стандарт определяет более миллиона кодовых позиций (от U+0000 до U+10FFFF).

Ключевые свойства Unicode:

- Первые 128 кодовых позиций (U+0000 – U+007F) полностью совпадают с **ASCII (American Standard Code for Information Interchange)**, обеспечивая частичную совместимость.
- Unicode сам по себе не является кодировкой. Он лишь определяет соответствие «символ ↔ число». Для представления этих чисел в виде байтов используются специальные кодировки (encodings).
- Стандарт имеет свои сложности: один и тот же видимый символ (глиф) может быть представлен либо одним кодпоинтом, либо комбинацией из базового символа и модифицирующего знака (например, диакритики). Например, символ 'ё' можно представить как U+0451 или как комбинацию 'е' (U+0435) и диерезиса (U+0308).

Кодировка UTF-8

Существует несколько способов кодирования кодпоинтов Unicode в байты: UCS-2, UCS-4, UTF-16. Наиболее популярным и де-факто стандартом в вебе и современных ОС стал **UTF-8 (Unicode Transformation Format, 8-bit)**.

Определение: UTF-8

UTF-8 (Unicode Transformation Format, 8-bit) — это кодировка Unicode с переменной длиной символа. Она кодирует кодпоинты в последовательности от 1 до 4 байт.

Основные преимущества UTF-8:

- **Совместимость с ASCII:** Любой текст в кодировке ASCII является корректным текстом в UTF-8, так как кодпоинты до U+007F кодируются одним байтом, полностью совпадающим с их ASCII-представлением.
- **Эффективность:** Маленькие значения кодпоинтов кодируются меньшим числом байт. Это экономит место для текстов на латинице.
- **Надёжность:** Нулевой байт (0x00) используется только для кодирования нулевого кодпоинта (U+0000). Это позволяет использовать стандартные C-функции для работы со строками, оканчивающимися нулём.
- **Самосинхронизация:** По значению любого байта можно определить, является ли он началом символа или частью многобайтовой последовательности. Это позволяет легко находить границы символов в потоке байт.

Принцип кодирования в UTF-8 основан на использовании старших битов первого байта для указания общей длины последовательности (см. [Table 5.2](#)).

5.2.4 Работа с Unicode в C/C++

Для поддержки Unicode в языках C и C++ существует тип `wchar_t` («широкий символ»).

Таблица 5.2 – Схема кодирования символов в UTF-8

Длина	Диапазон кодпоинтов	Схема байтов (x — биты кодпоинта)
1 байт	U+0000 – U+007F	0xxxxxxx
2 байта	U+0080 – U+07FF	110xxxxx 10xxxxxx
3 байта	U+0800 – U+FFFF	1110xxxx 10xxxxxx 10xxxxxx
4 байта	U+10000 – U+10FFFF	11110xxx 10xxxxxx 10xxxxxx 10xxxxxx

- `wchar_t`: Тип для хранения одного кодпоинта. Его размер зависит от платформы, но часто составляет 4 байта (32 бита), что достаточно для любого символа Unicode (аналогично кодировке UCS-4 в памяти).
- `L'a'`: Литерал типа `wchar_t`.
- `L"строка"`: Широкий строковый литерал (массив `const wchar_t`).
- Стандартная библиотека предоставляет аналоги для работы с широкими строками: `std::wstring`, `std::wcin`, `std::wcout` в C++ и функции `wprintf`, `wscanf` в C.

Локали

Чтобы стандартная библиотека знала, как преобразовывать широкие символы (`wchar_t`) в байтовые последовательности (например, в UTF-8) при вводе-выводе, ей необходимо сообщить текущие региональные настройки.

Определение: Локаль

локаль — это набор параметров, описывающих языковые и культурные особенности пользователя. **локаль** определяет кодировку текста (LC_TYPE), формат чисел, валюты, даты и времени. Она обычно задаётся через переменные окружения, такие как LANG или LC_ALL.

Перед началом работы с широким вводом-выводом в программе на C++ необходимо установить глобальную локаль, чтобы она была унаследована от настроек операционной системы.

```

1  #include <locale>
2  #include <iostream>
3  #include <string>
4
5  int main() {
6      // Set the global locale based on the system's environment variables.
7      // An empty string "" means "take from environment".
8      std::locale::global(std::locale(""));
9
10     // Now std::wcin and std::wcout will work correctly
11     // with the encoding specified in the locale (e.g., UTF-8).
12     std::wstring s;
13     std::wcout << L"input text: ";
14     std::wcin >> s;
15     std::wcout << L"you input: " << s << L", len: " << s.size() << L" symb" << std
        ::endl;
16
17     return 0;

```

18 }

Листинг 5.2 – Установка локали для корректной работы с Unicode в C++**Примечание**

Смешивать обычные потоки (`std::cout`) и широкие (`std::wcout`) в одной программе не рекомендуется. Такое смешивание может привести к непредсказуемому поведению и некорректному выводу, так как внутреннее состояние потоков может быть нарушено.

Итоги раздела

- Данные могут быть представлены в **текстовом** (человекочитаемом) или **бинарном** (машиночитаемом) формате.
- При работе с бинарными данными важно учитывать **порядок байтов** ([порядок байтов \(endianness\)](#)), особенно при передаче данных между разными системами.
- **ASCII** — исторически важная, но ограниченная 7-битная кодировка.
- **Unicode** является универсальным стандартом, присваивающим уникальный номер ([Unicode](#)) каждому символу.
- **UTF-8** — самая популярная кодировка для Unicode, эффективная и обратно совместимая с ASCII.
- В C/C++ для работы с Unicode используется тип `wchar_t` и связанные с ним строковые классы и функции ввода-вывода.
- Для корректного ввода-вывода Unicode-текста необходимо настроить **локаль** программы, чтобы она соответствовала системным настройкам.

Глава 6

6 Лекция

6.1 Представление целых чисел

В прошлых лекциях мы обсуждали представление текстовых данных. Теперь рассмотрим, как в памяти кодируются целые числа.

6.1.1 Беззнаковые числа

С беззнаковыми (unsigned) числами всё просто. Они представляются напрямую своим двоичным эквивалентом. Если у нас есть N бит, мы можем представить числа от 0 до $2^N - 1$. Например, для 3-битного числа:

- 000 \rightarrow 0
- 001 \rightarrow 1
- 010 \rightarrow 2
- 111 \rightarrow 7

6.1.2 Знаковые числа: Прямой код (Sign-Magnitude)

Первая и самая прямолинейная идея для представления знаковых чисел — использовать один бит (обычно старший) для кодирования знака, а остальные биты — для кодирования абсолютного значения (величины).

Например, для 3-битного числа (1 бит на знак, 2 на значение):

- 001 \rightarrow +1
- 010 \rightarrow +2
- 101 \rightarrow -1
- 110 \rightarrow -2

У этого подхода есть два существенных недостатка:

1. **Проблема двух нулей:** Существует два представления для нуля: 000 (+0) и 100 (-0). Это избыточно и усложняет проверки.
2. **Сложная арифметика:** Обычный двоичный сумматор "ломается". Сложение +1 (001) и -1 (101) в лоб даст 110, что равно -2, а не 0. Для выполнения арифметических операций требуются сложные проверки знаков.

6.1.3 Знаковые числа: Дополняющий код (Two's Complement)

Современные компьютеры решают эти проблемы, используя [Дополняющий код](#).

Определение: Дополняющий код

[Дополняющий код](#) — это способ представления знаковых чисел, основанный на арифметике по модулю 2^N , где N — количество бит.

- Положительные числа (и 0) представляются так же, как и беззнаковые (в диапазоне от 0 до $2^{N-1} - 1$).
- Отрицательные числа x (в диапазоне от -2^{N-1} до -1) представляются как беззнаковое число $2^N + x$.

Рассмотрим 3-битные числа (модуль $2^3 = 8$):

- 000 \rightarrow 0

- $001 \rightarrow 1$
- $010 \rightarrow 2$
- $011 \rightarrow 3$
- $100 \rightarrow 4$ (или $4 - 8 = -4$)
- $101 \rightarrow 5$ (или $5 - 8 = -3$)
- $110 \rightarrow 6$ (или $6 - 8 = -2$)
- $111 \rightarrow 7$ (или $7 - 8 = -1$)

Преимущества дополняющего кода:

- **Один ноль:** Значение 000 уникально.
- **Простая арифметика:** Обычный двоичный сумматор корректно работает как для знаковых, так и для беззнаковых чисел.

Примечание

Пример арифметики: Сложим 1 (001) и -2 (110) как знаковые.

$$001 + 110 = 111$$

Результат 111 в дополняющем коде — это -1 . Сложение работает. Теперь сложим 1 (001) и 6 (110) как беззнаковые.

$$001 + 110 = 111$$

Результат 111 в беззнаковом коде — это 7. Сложение также работает.

Получение отрицательного числа

Практическое правило для получения представления числа $-x$ из x в дополняющем коде:

1. Инвертировать все биты x (операция $\sim x$, побитовое НЕ).
2. Прибавить к результату 1.

Формула: $-x = \sim x + 1$.

Пример: Найти представление -3 (для 3-битного числа).

1. Берём 3: 011
2. Инвертируем (\sim): 100
3. Прибавляем 1: $100 + 1 = 101$

Результат 101 — это -3 , что совпадает с нашей таблицей.

6.2 Выравнивание данных в памяти

Определение: Выравнивание данных (Data Alignment)

Выравнивание данных — это ограничение, согласно которому данные определённого типа и размера должны размещаться в памяти по адресам, кратным некоторой степени двойки.

Например, 8-байтный `int64_t` должен иметь адрес, который делится на 8 (т.е. `address % 8 == 0`).

6.2.1 Зачем нужно выравнивание?

Выравнивание — это не просто прихоть компилятора, а требование, диктуемое аппаратным обеспечением (процессором).

- **Эффективность:** Процессоры читают данные из памяти не по одному байту, а "блоками" (например, по 4, 8 или 16 байт). Если 8-байтовое число "пересекает" границу такого блока (например, начинается с адреса 4 и заканчивается на 11), процессору придётся выполнить два чтения из памяти вместо одного.
- **Корректность:** На некоторых архитектурах (не x86) обращение по невыровненному адресу может привести к немедленному падению программы (аппаратному прерыванию). На x86 это "всего лишь" приводит к сильному замедлению.
- **Атомарность:** Операции чтения/записи по выровненным адресам, как правило, атомарны (неделимы). Невыровненная запись (например, 8 байт) может быть выполнена процессором как две отдельные записи по 4 байта.

Примечание

Проблема неатомарности особенно важна при работе с разделяемой памятью (shared memory). Представим, что два процесса (например, полученные через `fork()` с памятью `mmap(MAP_SHARED)`) работают с одним 8-байтным числом по невыровненному адресу.

Процесс А пишет новое значение. Он может успеть записать первые 4 байта, но не вторые. В этот момент Процесс Б читает это число и видит "мусор" — половину старого значения и половину нового.

Из-за этих требований компиляторы (C, C++, Rust, Go) автоматически вставляют "пропуски" (padding) в структуры, а стандартные аллокаторы (`malloc`, `operator new`) возвращают память, выровненную по максимальному требованию для стандартных типов (например, 16 байт на x86-64).

6.3 Процесс сборки программы

Рассмотрим, почему в C/C++ принято разделять код на заголовочные файлы (`.h`) и файлы реализации (`.cpp`).

6.3.1 Препроцессор и `#include`

Первый этап сборки — **Препроцессинг**. Директивы, начинающиеся с `#`, обрабатываются на этом этапе.

Директива `#include "header.h"` — это простая текстовая операция. Она заменяет эту строку содержимым файла `header.h`. Это можно проверить, запустив компилятор с флагом `-E`:

```
1 # gcc -E main.c
```

Листинг 6.1 – Запуск только препроцессора

Проблема многократного включения

Если один `.h` файл включается несколько раз (например, `a.h` и `b.h` оба включают `common.h`, а `main.cpp` включает `a.h` и `b.h`), мы получим дублирование кода и ошибки компиляции.

Для решения этой проблемы используются [Страж включения](#):

- **Классический способ (Стражи):**

```
1 #ifndef MY_HEADER_H
2 #define MY_HEADER_H
3
4 // ... soderzhimoe zagolovka ...
5
6 #endif // MY_HEADER_H
```

Листинг 6.2 – Использование `ifndef/define`

- **Современный способ:**

```
1 #pragma once
2
3 // ... soderzhimoe zagolovka ...
```

Листинг 6.3 – Использование `pragma once`

Оба способа гарантируют, что препроцессор включит тело файла только один раз.

6.3.2 Единицы трансляции и ускорение сборки

Основная причина разделения кода на `.h` и `.cpp` — это ****ускорение сборки**** больших проектов за счёт параллелизма.

Определение: Единица трансляции (Translation Unit)

Единица трансляции — это один `.c` или `.cpp` файл после того, как препроцессор "вклеил" в него содержимое всех `#include`.

Процесс сборки можно разбить на два этапа:

1. **Компиляция (Compilation):** Компилятор (например, `gcc -c`) *независимо и параллельно* обрабатывает каждую единицу трансляции, превращая её в **Объектный файл** (`<code>.o</code>`). Этот этап включает синтаксический анализ, оптимизации (`<code>-O2</code>`) и генерацию машинного кода.
2. **Линковка (Linking):** Линковщик (компоновщик) берёт все `<code>.o</code>` файлы и "сшивает" их в один исполняемый файл. Этот этап, как правило, последовательный, но он выполняется быстрее, чем полная перекомпиляция всего проекта.

Если мы меняем один `.cpp` файл, нам нужно перекомпилировать только его, а затем быстро перелинковать проект. Если бы весь код был в одном файле, любое изменение требовало бы полной перекомпиляции.

6.3.3 Объектные файлы и символы

Объектный файл (`.o`) — это "полуфабрикат". Он содержит машинный код, но в нём ещё нет информации о том, где находятся функции и переменные из *других* `.o` файлов.

Связь между файлами осуществляется через **символы**. С помощью утилиты `nm` можно посмотреть таблицу символов объектного файла.

```
1 # nm main.o
2 0000000000000000 T main
3                 U isEven
4                 U printf
5                 U scanf
```

Листинг 6.4 – Анализ символов с помощью `nm`

- **T (Text):** Символ *определён* (defined) в этом файле. Здесь определён `main`.
- **U (Undefined):** Символ *используется*, но не определён. Линковщик должен будет найти его в другом `.o` файле или библиотеке.

6.3.4 Линковка и релокации

Когда компилятор генерирует `main.o` и видит вызов `isEven()`, он не знает адреса этой функции. Вместо адреса он оставляет "дырку" — специальную запись, называемую **релокация**.

Задача линковщика:

1. Найти `main.o`, у которого `isEven` помечен как U.
2. Найти другой `.o` файл (например, `even.o`), у которого `isEven` помечен как T.
3. "Заполнить дырку" (выполнить релокацию) в `main.o`, подставив реальный адрес `isEven` из `even.o`.

Если линковщик не может найти определение для U-символа (или находит *несколько* определений), он выдаёт ошибку ("Undefined reference" или "Multiple definition").

Примечание

Поскольку `.o` файлы содержат только машинный код и таблицу символов (а не C++ или C код), они языково-независимы. Это позволяет компоновать программу из частей, написанных на разных языках (например, скомпилировать функцию на Rust, а вызвать её из C).

6.3.5 Формат ELF

В Linux исполняемые файлы и объектные файлы хранятся в формате **ELF (Executable and Linkable Format)**. Он состоит из **секций**, которые сообщают загрузчику ОС, как создать образ процесса в памяти.

Основные секции:

- **.text:** Исполняемый код (инструкции **CPU**). Загружается с правами "чтение + исполнение".

- **.rodata:** (Read-Only Data) Константные данные, например, строковые литералы (`"Hello"`). Загружается с правами "только чтение".
- **.data:** Инициализированные глобальные и статические переменные (`int x = 10;`). Загружается с правами "чтение + запись".
- **.bss:** Неинициализированные глобальные и статические переменные (`int y;`). Эта секция *не занимает места в файле*, она просто говорит загрузчику: "выдели X байт памяти и заполни их нулями".

6.4 Особенности C++: Имена и переменные

6.4.1 Искажение имён (Name Mangling)

В C++ можно объявлять функции с одинаковыми именами, но разными аргументами (перегрузка) или в разных пространствах имён:

```
1 void f();
2 void f(int);
3 namespace A { void f(); }
```

Линковщик C не справился бы с этим, так как он видит только один символ `f`. Компилятор C++ решает эту проблему, кодируя полную сигнатуру функции в имя символа. Этот процесс называется **Искажение имён (Name Mangling)**.

Например, `A::f()` может превратиться в `_Z1A1fv`.

6.4.2 Extern C

Чтобы C++ мог вызывать функции из C (или из ассемблера, который следует C-соглашениям) или наоборот, нужно отключить **Искажение имён (Name Mangling)**. Для этого используется `extern "C"`:

```
1 // Объявляем, что эта функция исполняет C ABI
2 // (без искажения имен)
3 extern "C" void my_c_function(int x);
```

6.4.3 Глобальные переменные: extern против static

Как и в случае с функциями, глобальные переменные нужно *объявлять* (в `.h`) и *определять* (в `.cpp`).

- **Правильный способ (Общая переменная):**

```
1 // Говорим компилятору, что переменная *где-то* существует
2 extern int shared_value;
```

```
1 // Выделяем память и задаем значение
2 int shared_value = 123;
```

Все `.cpp` файлы, включившие `def.h`, будут ссылаться на *одну и ту же* копию `shared_value`.

- **Неправильный способ (Локальные копии)**

```
1 // 'static' в global'ной области видимости
2 // делает переменную локальной для единицы трансляции
3 static int value = 0;
```

Если `main.cpp` и `other.cpp` включают `static.h`, *каждый* из них получит свою *собственную, независимую* копию `value`. Линкер не выдаст ошибки, но программа будет работать некорректно.

6.5 Основы ассемблера и архитектуры

6.5.1 Архитектура фон Неймана

Современные процессоры в основном следуют [Архитектура фон Неймана](#).

Ключевая особенность — единый блок памяти, в котором хранятся и данные, и инструкции (код) программы. Процессор выполняет инструкции последовательно, используя [Instruction Pointer, регистр-указатель на следующую инструкцию \(RIP\)](#) (Instruction Pointer) для отслеживания адреса текущей инструкции.

6.5.2 Регистры и память

Доступ к оперативной памяти (RAM) — медленная операция (порядка 100 ns). Чтобы процессор не простаивал, он содержит [Регистр](#) — сверхбыстрые ячейки памяти.

В архитектуре x86-64 (которую мы используем) есть 16 64-битных регистров общего назначения (`RAX`, `RCX`, `RDY`, `RSI`, `RDI`, `R8`...`R15` и т.д.).

Два регистра имеют особо важное значение:

- **RIP**: Указатель на инструкцию.
- **Stack Pointer, регистр-указатель на вершину стека (RSP)**: Указатель на вершину стека.

6.5.3 Стек и вызовы функций

Для реализации вызовов функций используется стек.

1. `call f` (Вызов функции):

- Процессор помещает адрес *следующей* за `call` инструкции (адрес возврата) на вершину стека (`push return_addr`).
- Процессор совершает безусловный переход на адрес функции `f` (`jmp f`).

2. `ret` (Возврат из функции):

- Процессор снимает адрес возврата с вершины стека (`pop return_addr`).
- Процессор совершает безусловный переход на этот адрес (`jmp return_addr`).

6.5.4 Соглашение о вызовах (ABI)

Как функции передают аргументы и возвращают значения? Процессор об этом "не знает". Это определяется программным соглашением — [ABI \(Application Binary Interface\)](#).

Для Linux x86-64 (System V ABI) действуют следующие правила:

Определение: Соглашение о вызовах (x86-64 System V ABI)

- Передача аргументов (целочисленных):
 - 1-й аргумент: `RDI`
 - 2-й аргумент: `RSI`

- 3-й аргумент: `<code>RDX</code>`
- 4-й аргумент: `<code>RCX</code>`
- 5-й аргумент: `<code>R8</code>`
- 6-й аргумент: `<code>R9</code>`
- **Передача аргументов (7-й и далее):**
 - Передаются через стек. Вызывающая сторона (caller) кладёт их на стек в обратном порядке до выполнения инструкции `call`.
- **Возвращаемое значение:**
 - `RAX`

Примечание

Аргументы на стеке. Поскольку `call` кладёт на стек адрес возврата, внутри вызываемой функции (callee) аргументы, переданные через стек, оказываются смещены:

- `[rsp]` — Адрес возврата (положен инструкцией `call`)
- `[rsp+8]` — 7-й аргумент
- `[rsp+16]` — 8-й аргумент
- и т.д.

(Здесь `[addr]` означает "прочитать 8 байт из памяти по адресу `addr`").

6.6 Практика: написание функций на ассемблере

Мы будем использовать синтаксис Intel. Файл `.S` должен начинаться с директив:

```
1 .intel_syntax noprefix # Ustanavlivaem sintaksis
2 .text # Nachalo sektsii koda
```

Чтобы сделать функцию `my_func` видимой для линковщика (C/C++), её нужно объявить глобальной:

```
1 .global my_func
2 my_func:
3     # ... instruktsii ...
4     ret
```

6.6.1 Пример 1: Возврат константы

```
1 // C++: extern "C" long constant();

3 .global constant
4 constant:
5     mov rax, 42 # Vozvrashchaemoe znachenie - v RAX
6     ret
```

6.6.2 Пример 2: Identity (аргумент -> возврат)

```
1 // C++: extern "C" long identity(long x);
```

```

7  .global identity
8  identity:
9      # 1-y argument 'x' prihodit v RDI
10     mov rax, rdi # Peremeshchaem RDI v RAX
11     ret

```

6.6.3 Пример 3: Сложение (два аргумента)

```

1 // C++: extern "C" long add(long x, long y);

12 .global add
13 add:
14     # x v RDI, y v RSI
15     add rdi, rsi # Skladyvaem: RDI = RDI + RSI
16     mov rax, rdi # Peremeshchaem rezul'tat (v RDI) v RAX
17     ret

```

6.6.4 Пример 4: Условный переход (If/Else)

```

1 /* C++: extern "C" long select(long cond, long a, long b);
2  * if (cond == 0) return b;
3  * else return a;
4  */

18 .global select
19 select:
20     # cond v RDI, a v RSI, b v RDX
21
22     # Proveryaem RDI na nol'.
23     # Operatsiya 'add' menyaet RFLAGS, v t.ch. Zero Flag (ZF)
24     add rdi, 0
25
26     # jz (Jump if Zero) - perehod, esli ZF=1 (rezul'tat byl 0)
27     jz .L_return_b
28
29 .L_return_a:
30     # cond != 0
31     mov rax, rsi # return a
32     ret
33
34 .L_return_b:
35     # cond == 0
36     mov rax, rdx # return b
37     ret

```

6.6.5 Пример 5: Цикл (Sum)

```

1 /* C++: extern "C" long sum(long n);
2  * long s = 0;
3  * for (long i = n; i > 0; i--) { s += i; }
4  * return s;
5  */

```

```
38 .global sum
39 sum:
40     # n в RDI
41     xor rax, rax # rax (summa) = 0
42
43 .L_loop_start:
44     add rax, rdi # summa += n
45     add rdi, -1 # n--
46
47     # 'add rdi, -1' (n--):
48     # - Esli n > 0, perenosa (carry) ne budet.
49     # - Esli n == 0, to 0 + (-1) daet perenos (borrow).
50
51     # jnc (Jump if No Carry) - prygnut', esli n > 0
52     jnc .L_loop_start
53
54     # n == 0, tsikl zavershen
55     ret
```

Итоги раздела

Итоги раздела "Ассемблер":

- Код на ассемблере — это прямое представление машинных инструкций (мнемоники).
- Взаимодействие с C/C++ происходит через [ABI \(Application Binary Interface\)](#) (соглашение о вызовах).
- В Linux x86-64 аргументы передаются через регистры (RDI, RSI и т.д.), а возвращаемое значение — через RAX.
- Аргументы, не поместившиеся в регистры (7-й и далее), передаются через стек и доступны по адресу [rsp+8], [rsp+16] и т.д.
- Управление потоком (if, loop) реализуется через [RFLAGS](#) и инструкции условных переходов (jz, jnc и др.).

Глава 7

7 Лекция

7.1 Адресация памяти в x86-64

Продолжаем изучение [низкоуровневый язык программирования, близкий к машинному коду \(Ассемблер\)](#). Ключевой темой является работа с памятью. В прошлый раз мы установили, что для обращения к памяти (разыменования) используется синтаксис с квадратными скобками.

7.1.1 Синтаксис Scale-Index-Base (SIB)

Общий синтаксис адресации памяти в 64-битном режиме (в [Intel-синтаксис](#)) выглядит следующим образом:

$$[rbase + rindex \times scale + displacement]$$

где:

- **rbase** — базовый регистр.
- **rindex** — регистр-индекс.
- **scale** — множитель (масштаб) для индекса. Допустимые значения: $scale \in \{1, 2, 4, 8\}$.
- **displacement** — константное смещение (сдвиг).

Этот синтаксис был разработан для удобной работы с массивами и структурами. Например, **rbase** может хранить адрес начала массива, **rindex** — индекс элемента, **scale** — размер одного элемента (e.g., 8 байт для `uint64_t`), а **displacement** — сдвиг до нужного поля внутри структуры.

```

1 ; * (uint64_t*)(rax + 8 * rdx) = rcx
2 ; (rax = base, rdx = index, 8 = scale)
3 mov [rax + rdx * 8], rcx
4
5 ; * (uint64_t*)(rbx + rbp + 32) = rax
6 ; (rbx = base, rbp = index, 1 = scale (default), 32 = displacement)
7 mov [rbx + rbp + 32], rax

```

Листинг 7.1 – Примеры SIB-адресации

7.1.2 Указание размера операнда

В листинг 7.1 ассемблер мог угадать размер операции (64 бита) по размеру регистра `rcx` или `rax`. Однако при работе с константами возникает неоднозначность.

```

1 mov [rax], 0 ; OSHIBKA: Neizvesten razmer: 1, 2, 4 ili 8 bayt?

```

Листинг 7.2 – Неоднозначность размера

Компилятор ассемблера не знает, какой размер данных вы намереваетесь записать. Для явного указания размера используются специальные директивы:

- **BYTE PTR** — 8 бит (1 байт).
- **WORD PTR** — 16 бит (2 байта).
- **DWORD PTR** — 32 бита (4 байта).
- **QWORD PTR** — 64 бита (8 байт).

```

1 ; * (uint32_t*)rax = 0
2 mov DWORD PTR [rax], 0

```

Листинг 7.3 – Явное указание размера (32 бита)

Итоги раздела

- Адресация SIB ($[base + index \times scale + disp]$) — основной механизм доступа к памяти.
- *scale* ограничен значениями {1, 2, 4, 8}.
- При неоднозначности (например, при записи константы) размер операции нужно указывать явно (e.g., DWORD PTR).

7.2 Инструкция LEA (Load Effective Address)

Инструкция [Load Effective Address](#), инструкция загрузки вычисленного адреса (LEA) — один из самых полезных и часто используемых инструментов в [Ассемблер](#).

Определение: LEA (Load Effective Address)

Инструкция `lea` **вычисляет** адрес, используя синтаксис SIB, но **не разыменовывает** его. Вместо этого она записывает вычисленный адрес в регистр-приемник.

```

1 ; MOV: Prochitat' 8 bayt po adresu [rax] i polozhit' v rdx
2 ; rdx = * (uint64_t*)rax
3 mov rdx, [rax]
4
5 ; LEA: Vychislit' adres (v etom sluchae prosto rax) i polozhit' v rdx
6 ; rdx = rax
7 lea rdx, [rax]

```

Листинг 7.4 – Сравнение MOV и LEA

Основное применение [LEA](#) — это вычисление адресов, но благодаря своей способности выполнять сложение и умножение (на 1, 2, 4, 8), она стала мощным инструментом для арифметических вычислений.

```

1 ; rdx = rax + 4 * rbx + 16
2 lea rdx, [rax + rbx * 4 + 0x10]

```

Листинг 7.5 – LEA для вычисления адреса

7.2.1 LEA как оптимизация компилятора

Компиляторы часто используют [LEA](#) для выполнения простых арифметических операций, так как [LEA](#) часто выполняется быстрее, чем инструкции умножения (такие как `imul`). Например, для компиляции функции `a * 3`:

```

1 uint64_t Mul3(uint64_t a) {
2     return a * 3;
3 }

```

Листинг 7.6 – C++ код для умножения на 3

Компилятор (g++ -O2) сгенерирует следующий код (листинг 7.7), используя **LEA** вместо умножения. В Linux (System V AMD64 ABI) первый аргумент (**a**) передается в регистре **rdi**, а возвращаемое значение — в **rax**.

$$a \times 3 = a \times (1 + 2) = a + a \times 2$$

Этот паттерн идеально ложится в SIB-адресацию: $[rdi + rdi \times 2]$.

```

1 0000000000000000 <Mul3(unsigned long)>:
2   0: f3 0f 1e fa endbr64 ; Zashchitnaya instruktsiya
3   4: 48 8d 04 7f lea rax,[rdi+rdi*2] ; rax = rdi + rdi * 2
4   8: c3 ret

```

Листинг 7.7 – Результат компиляции Mul3 (objdump)

Итоги раздела

- **lea** вычисляет адрес, но не читает память.
- Это мощный инструмент для компактных арифметических вычислений, часто используемый компиляторами.

7.3 Работа со стеком и локальными переменными

У процессора ограниченное количество регистров. При вызове функции (инструкция **call**) возникает проблема: как сохранить значения локальных переменных, если вызываемая функция может перезаписать ("испортить") регистры?

7.3.1 Проблема: Callee-clobbered регистры

Согласно соглашениям о вызовах (Calling Conventions), большинство регистров (как **rax**, **rdx**, **rdi** и т.д.) являются *callee-clobbered* — вызываемая функция (*callee*) имеет право изменять их без восстановления.

Рассмотрим код, где мы храним **a** и **b** в регистрах:

```

1 mov rax, 1 ; a = 1
2 mov rdx, 2 ; b = 2
3 call f ; f()
4 add rax, rdx ; ??? (rdx mozhet byt' isporchen funktsiey f)
5 ret

```

Листинг 7.8 – Проблема сохранения локальных переменных

После возврата из **f**, мы не можем полагаться на то, что в **rdx** все еще лежит 2.

7.3.2 Решение 1: Сохранение на стеке

Основной механизм для сохранения локальных переменных — это **стековый фрейм**. Мы можем "зарезервировать" место на стеке, сдвинув указатель стека **RSP**, и сохранить туда наши значения.

```

1 mov rax, 1 ; a = 1
2 mov rdx, 2 ; b = 2
3
4 sub rsp, 16 ; Rezerviruem 16 bayt na steke
5 mov [rsp + 8], rdx ; Sokhranyaem b (po smeshcheniyu 8)
6 mov [rsp], rax ; Sokhranyaem a (na vershinu steke)

```

```

7
8  call f ; f()
9
10 ; VOSSTANOVLENIE
11 mov rax, [rsp] ; Vosstanavlivaem a
12 mov rdx, [rsp + 8] ; Vosstanavlivaem b
13 add rsp, 16 ; Osvobozhdaem mesto na steke
14
15 add rax, rdx ; Teper' bezopasno
16 ret

```

Листинг 7.9 – Использование стека для локальных переменных

7.3.3 Решение 2: Callee-saved регистры

Некоторые регистры, напротив, являются *callee-saved* (например, `rbx`, `rbp`). Это означает, что если вызываемая функция хочет их использовать, она *обязана* сохранить их значение (обычно на стеке) и восстановить перед выходом (`ret`). Компиляторы используют это для оптимизации.

Рассмотрим C++ код:

```

1  uint64_t f();
2  uint64_t g();
3  uint64_t Sum() {
4      return f() + g();
5  }

```

Листинг 7.10 – C++ код Sum()

Чтобы вычислить `g()`, нужно сначала вызвать `f()`, но результат `f()` (который вернется в `rax`) будет перезаписан результатом `g()`. Компилятор (листинг 7.11) решает эту проблему, сохраняя результат `f()` в *callee-saved* регистре `rbx`.

```

1  <Sum(>:
2      push rbx ; 1. Sokhranit' staroe znachenie rbx
3      call <f(> ; 2. Vyzvat' f(). Rezul'tat v rax
4      mov rbx, rax ; 3. Spryatat' rezul'tat f() v rbx
5      call <g(> ; 4. Vyzvat' g(). Rezul'tat v rax
6      add rax, rbx ; 5. rax = rax + rbx (rezul'tat g() + rezul'tat f())
7      pop rbx ; 6. Vosstanovit' staroe znachenie rbx
8      ret

```

Листинг 7.11 – Дизассемблированный код Sum() (g++ -O2)

Примечание

В листинг 7.11 мы видим `call` на адрес вроде `<Sum()+0ха>` (в реальном `objdump` это часто `call 0`). Это **релокация**. На этапе компиляции адрес функции `f` еще неизвестен. Компилятор оставляет "дырку" (часто 0), а компоновщик (`linker`) на финальном этапе сборки подставляет в это место реальный адрес функции.

7.4 Фреймовые указатели (Frame Pointers)

В листинг 7.9 мы вручную двигали `RSP` (`sub rsp, 16`) и обращались к переменным относительно `RSP` (`[rsp + 8]`). Это работает, но усложняет отладку и трассировку стека.

Определение: Фреймовый указатель (RBP)

Base Pointer, регистр-указатель на базу стекового фрейма (RBP) (Base Pointer) — это регистр, который по соглашению используется для хранения адреса *начала* текущего *стековый фрейм*. Это обеспечивает "стабильный" якорь для доступа к локальным переменным, даже если RSP постоянно движется (например, при push/pop).

Для использования RBP применяется стандартный **пролог** (в начале функции) и **эпилог** (в конце).

```

1 f:
2   ; --- PROLOG ---
3   push rbp ; 1. Sokhranit' RBP predydushchey funktsii na stek
4   mov rbp, rsp ; 2. Zapomnit' tekushchuyu vershinu steka kak bazu (nachalo)
5                   ; nashego freyma. Teper' RBP stabilen.
6
7   ; --- Telo funktsii ---
8   ; Mesto dlya lokal'nykh peremennykh vydelyaetsya zdes'
9   ; sub rsp, 32 ; (vydelit' 32 bayta)
10  ; Dostup k peremennym idet otnositel'no RBP:
11  ; mov [rbp - 8], rax
12
13  ; --- EPILOG ---
14  mov rsp, rbp ; 1. Osvobodit' vse lokal'nye peremennye, vernuv rsp k baze
15  pop rbp ; 2. Vosstanovit' RBP predydushchey funktsii
16  ret

```

Листинг 7.12 – Стандартный пролог и эпилог функции

7.4.1 Структура стека с RBP

Когда каждая функция использует этот пролог, значения RBP на стеке образуют **односвязный список**. Каждое сохраненное значение RBP указывает на RBP предыдущей (вызвавшей) функции.

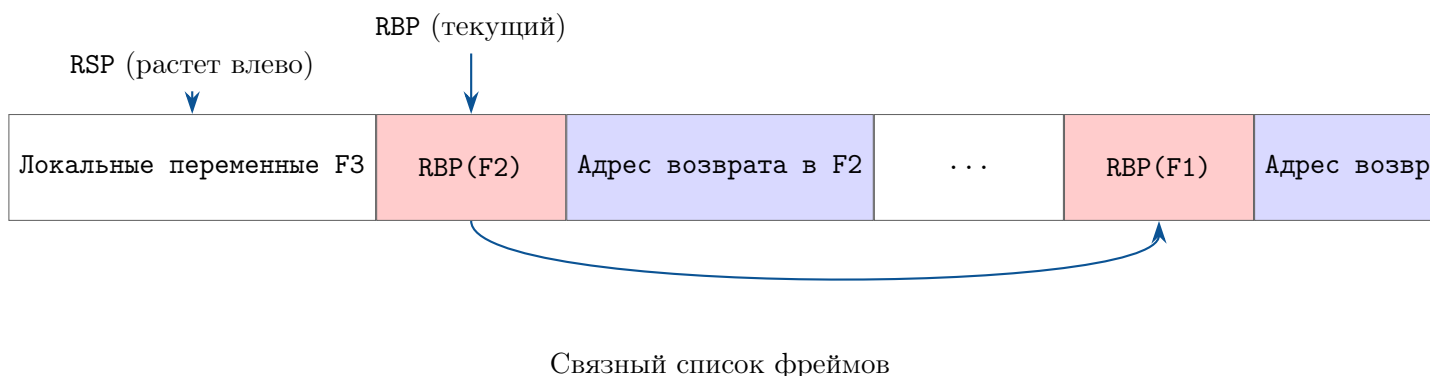


Рис. 7.1 – Структура стека при использовании фреймовых указателей (RBP)

Это позволяет отладчикам и другим инструментам легко "разматывать" стек (stack unwinding) и строить трассировку вызовов (stack trace).

Примечание

Использование **RBP** как фреймового указателя — это *соглашение*. Оно требует одного лишнего регистра и нескольких инструкций в прологе/эпилоге. Современные компиляторы (g++ -O2) по умолчанию часто отключают фреймовые указатели (-fomit-frame-pointer) для оптимизации. Вместо этого они генерируют специальную отладочную информацию (DWARF), которая позволяет разматывать стек, зная только **RIP**.

7.5 Секции данных в ассемблере

Ассемблерный код и данные не хранятся вперемешку. Они организованы в секции, которые сообщают операционной системе, как их следует загружать в память.

- **.text** — Код (инструкции). [cite: 216] Загружается с правами Read-Only и Execute (RX). [cite: 499]
- **.rodata** — Данные только для чтения. [cite: 217] (e.g., строковые литералы, константы). Загружаются с правами Read-Only (R). [cite: 501]
- **.data** — Инициализированные данные. [cite: 217] (e.g., глобальные переменные с начальным значением). Загружаются с правами Read-Write (RW). [cite: 505]
- **.bss** — Неинициализированные данные. [cite: 218] (e.g., `int x;`). Эти данные *не хранятся* в исполняемом файле, файл хранит только их размер. При загрузке ОС выделяет память и *обнуляет* ее. [cite: 507]

7.5.1 Директивы ассемблера для данных

Мы можем явно указать, в какую секцию помещать байты, с помощью директив.

```

1 ; Ob"yavlyаем sektsiyu .rodata
2 .section .rodata
3
4 .global s1 ; Delaem metku s1 vidimoy dlya linkera
5 s1:
6     .byte 0x48, 0x65 ; 'H', 'e'
7     .ascii "ll" ; 'l', 'l'
8     .asciz "o!" ; 'o', '!', i nulevoy bayt (terminator)

```

Листинг 7.13 – Пример секции .rodata (строка "Hello!")

В листинг 7.13 метка `s1` указывает на строку "Hello!\0".

Другие директивы для инициализации данных в секциях **.data** или **.rodata**:

```

1 .data
2     val1: .byte 0x10 ; 1 bayt
3     val2: .short 0x1234 ; 2 bayta
4     val3: .long 0x12345678 ; 4 bayta
5     val4: .quad 0x1122334455667788 ; 8 bayt
6
7     ; Zapolnit' 32 bayta znacheniem 0xFF
8     buffer: .space 32, 0xFF
9
10 .bss
11     ; Zarezervirovat' 1024 bayta (budut obnulyeny)

```

```
12 big_buffer: .skip 1024
```

Листинг 7.14 – Директивы .data и .bss

Примечание

Все есть байты. В конечном счете, ассемблер просто транслирует мнемоники инструкций в байты в секции `.text`. Можно написать функцию, используя только директиву `.byte`, если знать машинные коды.

7.6 Флаги процессора и условные переходы

Большинство арифметических и логических инструкций (ALU) изменяют специальный регистр флагов (EFLAGS/RFLAGS). Условные переходы (jcc) анализируют эти флаги.

Основные флаги, интересующие нас:

- **Carry Flag, флаг переноса (беззнаковое переполнение) (CF)** (Carry Flag) — Установлен, если произошел перенос/заём из старшего бита (индикатор **беззнакового** переполнения). [cite: 313, 545]
- **Zero Flag, флаг нуля (результат равен нулю) (ZF)** (Zero Flag) — Установлен, если результат операции равен нулю. [cite: 314, 546]
- **Sign Flag, флаг знака (установлен старший бит результата) (SF)** (Sign Flag) — Установлен, если старший бит результата равен 1 (индикатор **отрицательного** числа в знаковой интерпретации). [cite: 315, 546]
- **Overflow Flag, флаг переполнения (знаковое переполнение) (OF)** (Overflow Flag) — Установлен, если произошло **знаковое** переполнение (e.g., 100+100 дало отрицательный результат в 8-битном знаковом представлении). [cite: 316, 547]

```
1  # 0b0001 + 0b0111 = 0b1000 (1 + 7 = 8)
2  # Rezul'tat (8) imeet bit znaka (SF=1).
3  # Proizoshlo znakovoe perepolnenie (1+7 != -8) (OF=1).
4  # Flagi: SF, OF
5
6  # 0b1001 + 0b0111 = 0b0000 (s perenosom) (9 + 7 = 16)
7  # Rezul'tat 0 (ZF=1).
8  # Proizoshlo bezznakovoe perepolnenie (CF=1).
9  # Flagi: ZF, CF
10
11 # 0b0010 - 0b0011 = 0b1111 (s zaetom) (2 - 3 = -1)
12 # Rezul'tat -1 (SF=1).
13 # Proizoshel bezznakovyy zaem (CF=1).
14 # Flagi: CF, SF
```

Листинг 7.15 – Примеры установки флагов (4-битная арифметика)

Инструкция `cmp` (compare) — это, по сути, `sub`, которая не сохраняет результат, а только устанавливает флаги.

После `cmp` (или `add`, `sub`, `and`...) используются инструкции условного перехода jcc:

- `je / jz` — Jump if Equal / Jump if Zero (проверяет **ZF=1**).
- `jne / jnz` — Jump if Not Equal / Jump if Not Zero (**ZF=0**).

- js — Jump if Sign ($SF=1$).
- ja — Jump if Above (беззнаковое "больше") (проверяет $CF=0$ и $ZF=0$).
- jg — Jump if Greater (знаковое "больше") (проверяет $SF=OF$ и $ZF=0$).

7.7 Взаимодействие ассемблера и C/C++

Можно смешивать код на C/C++ и Ассемблер в одной программе, если соблюдать соглашения.

7.7.1 Позиционно-независимый код (PIC) и RIP-адресация

При попытке получить доступ к глобальной переменной (е.g., из C++ или секции `.data`) возникает проблема:

```

1 .data
2 my_var: .quad 123
3
4 .text
5 ; OSHIBKA: Ne budet rabotat' v sovremennykh OS
6 mov rax, [my_var]
```

Листинг 7.16 – Наивный доступ к глобальной переменной

Проблема в том, что в современных ОС из соображений безопасности (ASLR — Address Space Layout Randomization) программа загружается в память по *случайному* адресу. [cite: 587-589] Мы не знаем абсолютный адрес `my_var` на этапе компиляции. [cite: 583]

Решение — [Position-Independent Code](#), **позиционно-независимый код (PIC)** (PIC). Код не должен полагаться на абсолютные адреса, а только на *относительные*.

Определение: RIP-относительная адресация

В 64-битном режиме можно адресовать данные *относительно указателя инструкции (RIP)*. Так как **RIP** всегда указывает на следующую исполняемую инструкцию, а `my_var` находится на *неизменном* расстоянии от этой инструкции (весь код и данные сдвигаются вместе), этот сдвиг остается константой. [cite: 593-594]

```

1 extern c ; Ob"yavlyаем metku 'c' vneshney (opredelena v C++)
2
3 .text
4 GetC:
5     ; Korrektno: zagruzit' znachenie po adresu [rip + smeshchenie do 'c']
6     mov rax, [c+rip]
7     ret
```

Листинг 7.17 – Корректный доступ к глобальной переменной (PIC)

7.7.2 Оптимизация хвостового вызова (TCO)

Рассмотрим функцию-обертку, которая просто вызывает другую функцию и немедленно возвращает ее результат.

```

1 MyFuncWrapper:
2     ; ... podgotovka argumentov ...
3     call OtherFunc ; 1. Zapisat' adres vozvrata (A) na stek
4     ret ; 2. Snyat' adres (A) so steka i pereyti na nego
```

Листинг 7.18 – Неоптимальный хвостовой вызов

Здесь `call` кладет на стек адрес возврата (в `MyFuncWrapper`), а `ret` немедленно его снимает. Это лишняя работа.

[Tail Call Optimization](#), оптимизация хвостового вызова (TCO) (TCO) — это замена `call + ret` на один `jmp`.

```

1 MyFuncWrapper:
2     ; ... podgotovka argumentov ...
3     jmp OtherFunc ; Peredat' upravlenie OtherFunc

```

Листинг 7.19 – Оптимизированный хвостовой вызов (TCO)

Когда `OtherFunc` выполнит `ret`, она вернет управление не в `MyFuncWrapper`, а тому, кто вызвал `MyFuncWrapper` (т.к. его адрес возврата все еще лежит на вершине стека). [cite: 608] Компиляторы (-O1 и выше) активно применяют эту оптимизацию.

7.7.3 Вызов функций C (`scanf` / `printf`)

Пользоваться вводом-выводом C++ (`iostream`) из [Ассемблер](#) почти невозможно из-за name mangling (искажения имен). [cite: 638-640] Гораздо проще использовать функции из C `<stdio.h>`, такие как `scanf` и `printf`. [cite: 643]

При этом нужно строго соблюдать два правила соглашения о вызовах (ABI): **1. Аргументы:** Первые 6 целочисленных аргументов/указателей передаются через регистры (именно в таком порядке): RDI, RSI, RDX, RCX, R8, R9.

2. Выравнивание стека: Перед инструкцией `call` [RSP](#) (указатель стека) должен быть выровнен по 16-байтной границе. [cite: 602]

Примечание

Ловушка выравнивания: Когда нашу функцию `main` вызывают, [RSP](#) уже выровнен по 16-байтной границе. Но инструкция `call` (которая вызвала `main`) помещает на стек 8-байтный адрес возврата. [cite: 603] Это означает, что *внутри* нашей функции `main` [RSP](#) не выровнен (он равен $16N + 8$). Перед тем, как мы сами сделаем `call` (например, `call scanf`), мы должны "скомпенсировать" эти 8 байт, например, `sub rsp, 8`. [cite: 604, 733]

```

1 .intel_syntax noprefix
2
3 .section .rodata
4 ; Formatnaya stroka dlya chteniya ("%lld")
5 read_fmt: .asciz "%lld"
6 ; Formatnaya stroka dlya zapisi ("%lld\n")
7 write_fmt: .asciz "%lld\n"
8
9 .text
10 .global main
11 main:
12     ; --- PROLOG ---
13     ; Vydelyaem 8 bayt dlya peremennoy 'n'
14     ; I zaoschno VYRAVNIVAEM stek (rsp byl 16N+8, stal 16N)
15     sub rsp, 8

```

```
16
17 ; --- Vyzov scanf ---
18 ; scanf("%lld", &n);
19 ; &n teper' = adres [rsp]
20
21 ; Arg 1 (RDI): Adres formatnoy stroki
22 lea rdi, [read_fmt+rip]
23 ; Arg 2 (RSI): Adres, kuda pisat' rezul'tat (vershina steka)
24 mov rsi, rsp
25
26 ; Dlya variadic funktsiy (kak scanf) nuzhno obnulit' rax
27 xor rax, rax
28 call scanf
29
30 ; --- Vyzov printf ---
31 ; printf("%lld\n", n + 1);
32 ; Zagruzhaem 'n' so steka
33 mov rsi, [rsp]
34 ; Uvelichivaem
35 add rsi, 1
36
37 ; Arg 1 (RDI): Adres formatnoy stroki
38 lea rdi, [write_fmt+rip]
39 ; Arg 2 (RSI): Znachenie (n + 1)
40 ; (uzhe v rsi)
41
42 xor rax, rax
43 call printf
44
45 ; --- EPILOG ---
46 ; "return 0;"
47 ; Po ABI, my vozvrashchaem znachenie iz main cherez RAX
48 xor rax, rax
49
50 ; Osvobozhdaem mesto na steke
51 add rsp, 8
52 ret
```

Листинг 7.20 – Пример: чтение числа (n) и вывод (n+1) на Assembler

Итоги раздела

- Для доступа к глобальным данным используйте **RIP-относительную адресацию** ([my_var+rip]).
- `call func + ret` можно заменить на `jmp func` (TCO).
- При вызове функций C (e.g., `printf`) стек должен быть выровнен по 16 байт до инструкции `call`.
- Аргументы передаются через RDI, RSI, RDX, RCX...
- `scanf` ожидает *указатель* (адрес) в RSI, `printf` — *значение*.
- Возвращаемое значение из `main` — это то, что лежит в RAX в момент `ret`.

7.8 Синтаксисы ассемблера: Intel vs. AT&T

Существует два доминирующих синтаксиса x86 [Ассемблер](#).

- **Intel-синтаксис:** (Используется в этой лекции, в документации Intel, Microsoft).
- **AT&T-синтаксис (GNU):** (Используется по умолчанию в `objdump` и `gcc`). [cite: 688]

Ключевые отличия:

Таблица 7.1 – Сравнение синтаксисов Intel и AT&T (GNU)

Аспект	Intel (мы)	AT&T (GNU)
Порядок операндов	<code>mov rax, rbx</code> (Приемник, Источник)	<code>mov %rbx, %rax</code> (Источник, Приемник)
Регистры	<code>rax, rbx</code>	<code>%rax, %rbx</code> (с префиксом %)
Константы	<code>16, 0x10</code>	<code>\$16, \$0x10</code> (с префиксом \$)
Адресация	<code>[rax + rbx * 4 + 32]</code>	<code>32(%rax, %rbx, 4)</code>
Размер	<code>DWORD PTR [rax]</code>	<code>movl \$0, (%rax)</code> (суффикс l/q/w/b)

Примечание

Полезно уметь читать оба синтаксиса. В `objdump` можно включить [Intel-синтаксис](#) с помощью флага `-M intel`.

Глава 8

8 Лекция

8.1 Оптимизация ассемблерного кода

Завершающая лекция по ассемблеру посвящена методам его эффективного использования, взаимодействию с ядром и компоновщиком, а также созданию программ, полностью независимых от стандартной библиотеки.

8.1.1 Проблема раздельной компиляции

Ранее мы рассматривали вызов ассемблерной функции из C++ с использованием раздельной компиляции. Этот подход имеет существенные недостатки производительности.

Рассмотрим пример с подсчётом суммы арифметической прогрессии. Если функция подсчёта реализована в C++ (и компилируется Clang), компилятор может распознать паттерн и заменить цикл на формулу $O(1)$. Если же функция вынесена в отдельный .S файл, компилятор видит только её объявление и вынужден генерировать цикл $O(n)$.

Более того, сам вызов функции между единицами трансляции (translation units) — это не бесплатная операция. Он включает:

- Инструкцию `call`, которая сохраняет адрес возврата в стек и выполняет переход (jump).
- Инструкцию `ret`, которая извлекает адрес из стека и выполняет косвенный переход ([косвенный переход](#)) по нему.

Эти операции вносят накладные расходы, которых можно избежать.

8.1.2 Встроенный ассемблер (GNU Inline Assembly)

Для устранения накладных расходов на вызов функции можно использовать [встроенный ассемблер](#). Эта конструкция позволяет компилятору вставить ассемблерный код непосредственно в тело C++ функции, избегая `call/ret`.

Определение: Синтаксис GNU Inline Assembly

Конструкция `asm` в C/C++ (GCC, Clang) имеет следующий расширенный синтаксис:

```
asm [volatile] ("assembly template"
                : output operands    /* optional outputs */
                : input operands     /* optional inputs */
                : clobber list       /* optional clobbers */
);
```

- **assembly template:** Строковый литерал с ассемблерным кодом. Входы и выходы подставляются как `%0`, `%1` и т.д.
- **output operands:** Список переменных C/C++, в которые нужно записать результат.
- **input operands:** Список переменных/выражений C/C++, которые нужно передать в ассемблер.
- **clobber list:** Список регистров или состояний, которые изменяются (портятся) внутри вставки, о чём компилятор должен знать.

Рассмотрим пример сложения двух чисел (листинг 8.1).

```
1 long add_asm(long a, long b) {
2     long res;
```

```

3  asm (
4      "mov %[a_reg], %[res_reg]\n\t" // res = a
5      "add %[b_reg], %[res_reg]\n\t" // res += b
6      : [res_reg] "=&r" (res) // Output: res, in any register (r)
7                          // & = early clobber
8      : [a_reg] "r" (a), // Input: a, in any register (r)
9        [b_reg] "r" (b) // Input: b, in any register (r)
10     : "cc" // Clobbers: "cc" (condition codes / flags)
11 );
12 return res;
13 }

```

Листинг 8.1 – Использование inline asm для сложения

Операнды и ограничения (Constraints)

Компилятор не понимает семантику ассемблерного кода; для него это просто шаблон. Мы должны явно описать интерфейс между C++ и ассемблером с помощью ограничений:

- **"r"**: Поместить переменную в регистр общего назначения (например, `eax`, `rdi`).
- **-r"**: Выходной операнд (=), который будет в регистре.
- **"&r"**: Ограничение **Early Clobber (&)**. Оно сообщает компилятору, что этот выходной регистр (`res_reg`) будет перезаписан *до* того, как все входные операнды (`a_reg`, `b_reg`) будут использованы. Это запрещает компилятору выделять один и тот же физический регистр для `res` и, например, `a`.

Список порчи (Clobbers)

Ассемблерная вставка может иметь побочные эффекты. Инструкция `add` изменяет регистр флагов (EFLAGS/RFLAGS). Если мы не сообщим об этом компилятору, он может ошибочно предположить, что флаги, установленные *до* `asm`-вставки, останутся неизменными *после* неё.

- **"cc"**: Сообщает компилятору, что регистр флагов (condition codes) был изменён.
- **"memory"**: Сообщает, что вставка читает или пишет в память по адресам, неизвестным компилятору. (См. раздел 8.2.3).
- **"rax", "rcx"** и т.д.: Сообщает, что конкретный регистр был изменён.

8.1.3 Оптимизация на этапе компоновки (LTO)

Встроенный ассемблер решает проблему вызова, но не проблему оптимизации. Если функция `add` находится в другом `.cpp` файле, компилятор всё ещё не видит её реализацию при компиляции `main.cpp` и не может, например, заинлайнить её.

Определение: Link Time Optimization (LTO)

Link Time Optimization (оптимизация на этапе компоновки) (LTO) — это техника, при которой компилятор генерирует объектные файлы не в виде машинного кода, а в виде **Intermediate Representation (промежуточное представление) (IR)**. На этапе компоновки (линковки) компилятор снова запускается, считывает **IR** из *всех* объектных файлов и выполняет оптимизации (включая инлайнинг, удаление мёртвого кода, константное сворачивание) так, как если бы весь код находился в одной единице трансляции.

Инфраструктура [Low Level Virtual Machine](#) (инфраструктура для построения компиляторов) ([LLVM](#)) (используемая Clang) идеально для этого подходит.

- **Без LTO:** `clang++ -O2 → main.o (x86-64), add.o (x86-64)`. Линкер просто склеивает их.
- **С LTO (-flto):** `clang++ -flto -O2 → main.o (LLVM IR), add.o (LLVM IR)`. На этапе линковки `clang` видит IR обеих функций, инлайнит `add` в `main` и может применить оптимизацию (например, свернуть сумму арифметической прогрессии в константу).

Примечание

Объектные файлы LTO, сгенерированные Clang, не являются стандартными ELF-файлами с машинным кодом. Это архивы [LLVM IR](#) (LLVM-AR). Для их просмотра вместо `objdump` используется `llvm-dis`. GCC также поддерживает LTO, но обычно встраивает своё IR (GIMPLE) в специальные секции ELF-файлов.

Итоги раздела

- Вызовы функций между `.cpp` и `.S` файлами несут накладные расходы (`call/ret`).
- [встроенный ассемблер](#) позволяет встроить ассемблерный код в C++, устраняя эти расходы, но требует аккуратного описания интерфейса (входы, выходы, [clobbers](#) (список порчи)).
- [LTO](#) позволяет компилятору оптимизировать код *между* единицами трансляции, генерируя промежуточное [IR](#) вместо машинного кода.

8.2 Взаимодействие с ядром и побочные эффекты

8.2.1 Прямой вызов syscall

Ассемблерные вставки позволяют нам выполнять **Системный вызов** напрямую, минуя обёртки стандартной библиотеки (libc).

Определение: Соглашение о syscall в Linux x86-64

- Инструкция: `syscall`.
- Номер системного вызова: передаётся в `RAX`.
- Аргументы (по порядку): `RDI`, `RSI`, `RDX`, `R10`, `R8`, `R9`.
- Возвращаемое значение: в `RAX`.
- **Порча:** Инструкция `syscall` уничтожает содержимое `RCX` и `R11`.

Примечание

Соглашение о `syscall` отличается от стандартного System V ABI в 4-м аргументе (`R10` вместо `RCX`). Это связано с тем, что `syscall` использует `RCX` для сохранения адреса возврата (`RIP`) и `R11` для сохранения `RFLAGS`, чтобы ядро могло вернуться в пользовательский процесс с помощью инструкции `sysret`.

В листинг 8.2 показан пример вызова `write` (номер 1) для печати "Hello".

```
1 #include <sys/syscall.h> // for SYS_write
2 #include <unistd.h> // for STDOUT_FILENO
3
4 long write_syscall(int fd, const char* buf, size_t count) {
5     long ret;
6     asm volatile (
7         "syscall"
8         : "=a" (ret) // Output: in RAX (a)
9         : "a" (SYS_write), // Input: syscall number in RAX (a)
10        "D" (fd), // Input: arg1 in RDI (D)
11        "S" (buf), // Input: arg2 in RSI (S)
12        "d" (count) // Input: arg3 in RDX (d)
13        : "rcx", "r11", "memory" // Clobbers: syscall clobbers rcx, r11
14                                   // "memory" because 'buf' is read
15    );
16    return ret;
17 }
18
19 int main() {
20     const char* msg = "Hello from syscall!\n";
21     write_syscall(STDOUT_FILENO, msg, 20);
22     return 0;
23 }
```

Листинг 8.2 – Системный вызов `write` через inline asm

8.2.2 volatile и побочные эффекты

Что произойдёт, если мы вызовем `write_syscall`, но не будем использовать возвращаемое значение (`ret`)?

```
1 int main() {
2     const char* msg = "Hello...\n";
3     // The compiler (with -O2) might DELETE this line!
4     write_syscall(STDOUT_FILENO, msg, 10);
5     return 0;
6 }
```

Листинг 8.3 – Проблема оптимизации

С точки зрения компилятора, функция `write_syscall` (листинг 8.2 без `volatile`) — это "чёрный ящик который принимает 4 аргумента и возвращает `long`. Если этот `long` не используется, компилятор вправе удалить вызов целиком, следуя правилу "as-if" (программа должна вести себя *так, как если бы* она выполнялась).

Проблема в том, что у `syscall` есть **побочный эффект** (вывод на экран), о котором компилятор не знает.

Определение: `asm volatile`

Ключевое слово `volatile` перед `asm` (`asm volatile (...)`) запрещает компилятору:

1. **Удалять** эту ассемблерную вставку, даже если её выходные операнды не используются.
2. **Переупорядочивать** её относительно других `volatile` операций (например, доступа к `volatile` переменным).

Это необходимо для всех ассемблерных вставок, имеющих побочные эффекты (side effects), такие как системные вызовы.

8.2.3 Использование `asm` для барьеров компиляции

Конструкцию `asm volatile` можно использовать для управления оптимизациями компилятора.

Барьер оптимизации (`DoNotOptimize`)

Иногда в бенчмарках нужно C++ значение, чтобы компилятор не "выкинул" всё вычисление этого значения.

```
1 // Helper function, similar to Google Benchmark
2 template <class T>
3 void DoNotOptimize(T const& value) {
4     // The asm block does nothing, but "reads" 'value'
5     // 'm' = memory operand
6     asm volatile("" :: "m" (value) : "memory");
7 }
8
9 // ...
10 long result = complex_calculation();
11 DoNotOptimize(result); // Now the compiler MUST
12                        // compute 'result'
```

Листинг 8.4 – Запрет оптимизации переменной

Барьер памяти (Compiler Fence)

```
asm volatile("" ::: "memory");
```

Листинг 8.5 – Барьер памяти компилятора

Список `clobbers` (список порчи), содержащий `"memory"`, сообщает компилятору, что эта вставка может читать или писать в *любую* ячейку памяти. Это заставляет компилятор:

- **Сбросить** (spill) все значения из регистров, которые были изменены, обратно в память *до* этой вставки.
- **Загрузить** (reload) значения из памяти *после* этой вставки, если они понадобятся, не полагаясь на кэшированные в регистрах значения.

Это барьер *только для компилятора*, он не генерирует инструкций барьера **CPU** (типа `mfence`).

Итоги раздела

- Системные вызовы в Linux x86-64 выполняются инструкцией `syscall`, используя регистры `RAX`, `RDI`, `RSI`, `RDY`, `R10`...
- Инструкция `syscall` портит `RCX` и `R11`.
- `asm volatile` необходимо использовать, когда вставка имеет побочные эффекты (как `syscall`), чтобы компилятор её не удалил.
- `asm volatile` с `memory` в `clobbers` служит барьером памяти для компилятора.

8.3 Указатели, функции и полиморфизм

Знание ассемблера позволяет понять, как реализованы высокоуровневые конструкции C++, такие как указатели на функции и виртуальные методы.

8.3.1 Указатели на функции и косвенные переходы

Указатель на функцию в C++ — это переменная, хранящая адрес.

```
1 int f1(int x) { return x; }
2 int f2(int x) { return x * 2; }
3 int f3(int x) { return x * 3; }
4
5 // Syntax: ret_type (*var_name)(arg_types)
6 int (*fn_array[])(int) = { f1, f2, f3 };
7
8 int main() {
9     int index = 1; // Assume this came from user input
10    // ...
11    int result = fn_array[index](5); // calls f2(5)
12 }
```

Листинг 8.6 – Массив указателей на функции

Во что транслируется `fn_array[index](5)`?

1. Загрузка адреса из `fn_array[index]` в регистр (например, `rax`).
2. Загрузка аргумента 5 в `rdi`.
3. Выполнение **косвенный переход**: `call rax`.

Определение: Указатель на функцию и косвенный переход

Численное значение указателя на функцию — это, как правило, адрес первой инструкции этой функции в секции `.text`. Вызов по такому указателю реализуется CPU через **косвенный вызов** (`call <reg>`), адрес которого неизвестен на этапе компиляции.

8.3.2 Защита от атак: `endbr64`

Косвенные переходы — основной вектор атак (ROP/JOP), когда злоумышленник получает контроль над регистром (`rax`) и заставляет программу прыгнуть не на начало функции, а в середину другой функции (на "гаджет").

Для борьбы с этим в современных CPU (Intel CET) введена инструкция `endbr64`.

- Компиляторы (GCC/Clang) теперь вставляют `endbr64` в начало каждой функции.
- Если ОС и CPU включают защиту, любой **косвенный переход** (`call rax`), который приземляется не на инструкцию `endbr64`, вызовет аппаратное исключение (fault).
- Это гарантирует, что косвенные вызовы могут приземляться только на легитимные начала функций.

Примечание

На данный момент (в лекции) Linux использует эту защиту в основном для кода ядра, но не для пользовательских (userspace) приложений. Однако компиляторы всё

равно генерируют `endbr64` для совместимости в будущем.

8.3.3 Реализация виртуальных функций C++

Динамический полиморфизм в C++ (ключевое слово `virtual`) также построен на косвенных вызовах.

Определение: `vptr` и `vtable`

- **`vtable` (Таблица виртуальных методов):** Статический массив указателей на функции, создаваемый компилятором для *каждого класса*, имеющего виртуальные методы.
- **`vptr` (Указатель на `vtable`):** Скрытый указатель, добавляемый компилятором в *каждый объект* такого класса. `vptr` указывает на `vtable`, соответствующую реальному типу объекта.

Рассмотрим вызов `a->foo()`:

```

1 struct A {
2     virtual void foo() { /* A's foo */ }
3 };
4 struct B : A {
5     virtual void foo() override { /* B's foo */ }
6 };
7
8 int main() {
9     A* a = new B();
10    a->foo(); // <-- How does this work?
11 }

```

Листинг 8.7 – Виртуальный вызов

Вызов `a->foo()` (где `a` в `rdi`) транслируется в:

1. `mov rax, [rdi]` ; Загрузить `vptr` из объекта (`this`) в `rax`
2. `call [rax]` ; Вызвать функцию по первому адресу в `vtable`

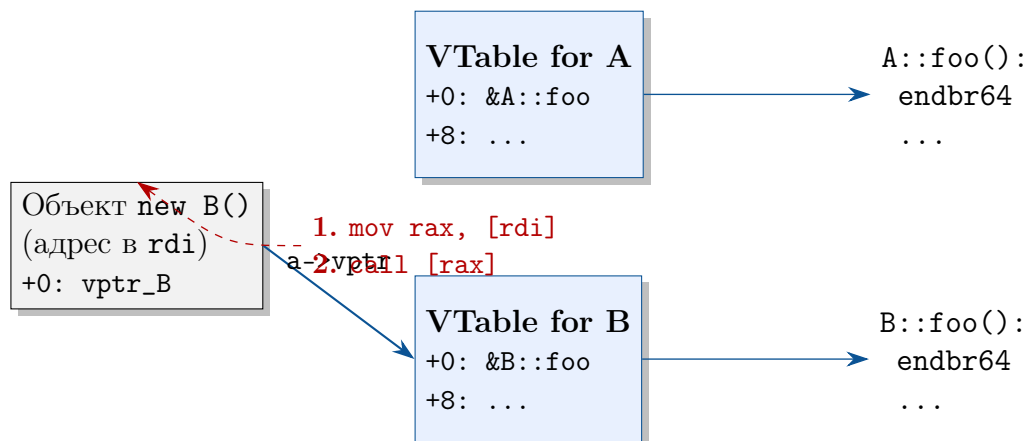


Рис. 8.1 – Схема виртуального вызова через `vptr` и `vtable`

Цена виртуализации

- **Цена по памяти:** +8 байт на *каждый* объект (`vptr`). Это нарушает принцип C++ "платишь только за то, что используешь т.к. вы платите за `vptr`, даже если никогда не делаете виртуальных вызовов.
- **Цена по времени:** Виртуальный вызов требует двух обращений к памяти (чтение `vptr`, чтение адреса из `vtable`) и *косвенный переход*, что медленнее прямого `call`.

8.3.4 JIT-компиляция (Just-in-Time)

Зная, что код — это просто байты в памяти, мы можем генерировать его во время выполнения.

Определение: JIT-компиляция

Just-in-Time (компиляция «на лету») (JIT) — это техника, при которой машинный код генерируется не на этапе компиляции, а во время выполнения программы. Это позволяет создавать код, оптимизированный под конкретные данные (например, SQL-запрос в ClickHouse) или под конкретное железо (например, используя AVX-инструкции, если они доступны на CPU пользователя).

Процесс JIT в Linux:

1. Выделить память с помощью `mmap` с правами `PROT_READ | PROT_WRITE`.
2. Записать в эту память байты машинного кода.
3. Изменить права памяти с помощью `mprotect` на `PROT_READ | PROT_EXEC (W⊕X)`.
4. Преобразовать указатель на эту память в указатель на функцию.
5. Вызвать сгенерированную функцию.

```
1 #include <sys/mman.h> // mmap, mprotect
2 #include <string.h> // memcpy
3
4 // Bytes for the function:
5 // mov rax, rdi (48 89 f8)
6 // add rax, rsi (48 01 f0)
7 // ret (c3)
8 unsigned char code[] = { 0x48, 0x89, 0xf8, 0x48, 0x01, 0xf0, 0xc3 };
9
10 typedef long (*add_func_t)(long, long);
11
12 int main() {
13     void* mem = mmap(NULL, sizeof(code),
14                     PROT_READ | PROT_WRITE,
15                     MAP_PRIVATE | MAP_ANONYMOUS, -1, 0);
16
17     memcpy(mem, code, sizeof(code));
18
19     // Important: make the memory executable
20     mprotect(mem, sizeof(code), PROT_READ | PROT_EXEC);
21
22     add_func_t fn = (add_func_t)mem;
23     long result = fn(10, 20); // result == 30
```

```
24  
25     munmap(mem, sizeof(code));  
26     return 0;  
27 }
```

Листинг 8.8 – Ручная JIT-компиляция функции add(a, b)

Примечание

Приложения вроде ClickHouse или JVM не пишут байты вручную. Они встраивают в себя бэкенд компилятора (например, [LLVM](#)) и используют его API для генерации оптимизированного кода "на лету".

Итоги раздела

- Указатели на функции реализуются через [косвенный переход](#) (call rax).
- Инструкция [endbr64](#) защищает от атак, пометая легитимные цели для таких переходов.
- Виртуальные вызовы C++ используют [vptr](#) (в объекте) и [vtable](#) (на класс) для реализации [косвенный переход](#), что несёт расходы памяти и времени.
- [JIT](#)-компиляция позволяет генерировать машинный код во время выполнения с помощью `mmap` и `mprotect`.

8.4 Динамическая компоновка

динамическая библиотека (.so) — это код, который компоуется с программой не при сборке, а при запуске.

8.4.1 Мотивация и основы (.so)

Две основные причины для использования **динамическая библиотека (.so)**:

1. **Экономия памяти:** Множество программ (bash, ls, g++) используют одну и ту же стандартную библиотеку (libc, libstdc++). Вместо того чтобы каждый процесс загружал свою копию, ОС загружает **.so** в память один раз и отображает её в адресные пространства всех процессов.
2. **Оптимизация под платформу:** Можно иметь несколько реализаций **тепсру** (обычную, SSE, AVX) и при запуске загрузчик выберет ту **.so**, которая оптимизирована под текущий **CPU**.

Загрузчик (ld.so в Linux) отвечает за поиск и загрузку всех зависимостей (их можно посмотреть командой `ldd a.out`) перед запуском `_start`.

Позиционно-независимый код (PIC)

Динамическая библиотека не знает, по какому адресу она будет загружена в виртуальную память. Поэтому её код не может использовать абсолютную адресацию.

Определение: Position Independent Code (PIC)

PIC — это код, который использует **относительную адресацию** (в x86-64 — относительно регистра RIP) для всех переходов и доступа к данным. Это позволяет загружать **.so** в любое место в памяти без необходимости её модификации. Для сборки **PIC** используется флаг `-fPIC`.

8.4.2 Механизмы PLT и GOT

Как `main` (скомпилированный) может вызвать `printf` (адрес которой станет известен только при запуске)?

Определение: GOT и PLT

- **Global Offset Table (глобальная таблица смещений) (GOT) (Global Offset Table):** Глобальная таблица смещений. Это массив в секции данных, хранящий *реальные адреса* внешних функций и переменных.
- **Procedure Linkage Table (таблица компоновки процедур) (PLT) (Procedure Linkage Table):** Таблица компоновки процедур. Это секция *исполняемого* кода, содержащая "трамплины" (stubs) — по одному на каждую внешнюю функцию.

По умолчанию в Linux используется **ленивое связывание** (ленивое связывание).

Процесс первого вызова printf:

1. `main` вызывает не `printf`, а трамплин `printf@plt`.
2. Трамплин `printf@plt` прыгает на адрес, указанный в **GOT** для `printf`.
3. *Изначально* **GOT** указывает не на `printf`, а обратно на код в **PLT**.

4. Этот код в **PLT** кладёт ID функции (`printf`) в стек и прыгает на **динамический резолвер** (часть `ld.so`).
5. Резолвер находит реальный адрес `printf` в загруженной `libc.so`.
6. (**Патчинг**) Резолвер *перезаписывает* запись `printf` в **GOT**, указывая на реальный адрес.
7. Резолвер прыгает на реальный `printf`.

Второй и последующие вызовы `printf`:

1. `main` вызывает `printf@plt`.
2. Трамплин `printf@plt` прыгает на адрес, указанный в **GOT**.
3. **GOT** уже содержит реальный адрес `printf`. Происходит прямой переход к `printf`, минуя резолвер.

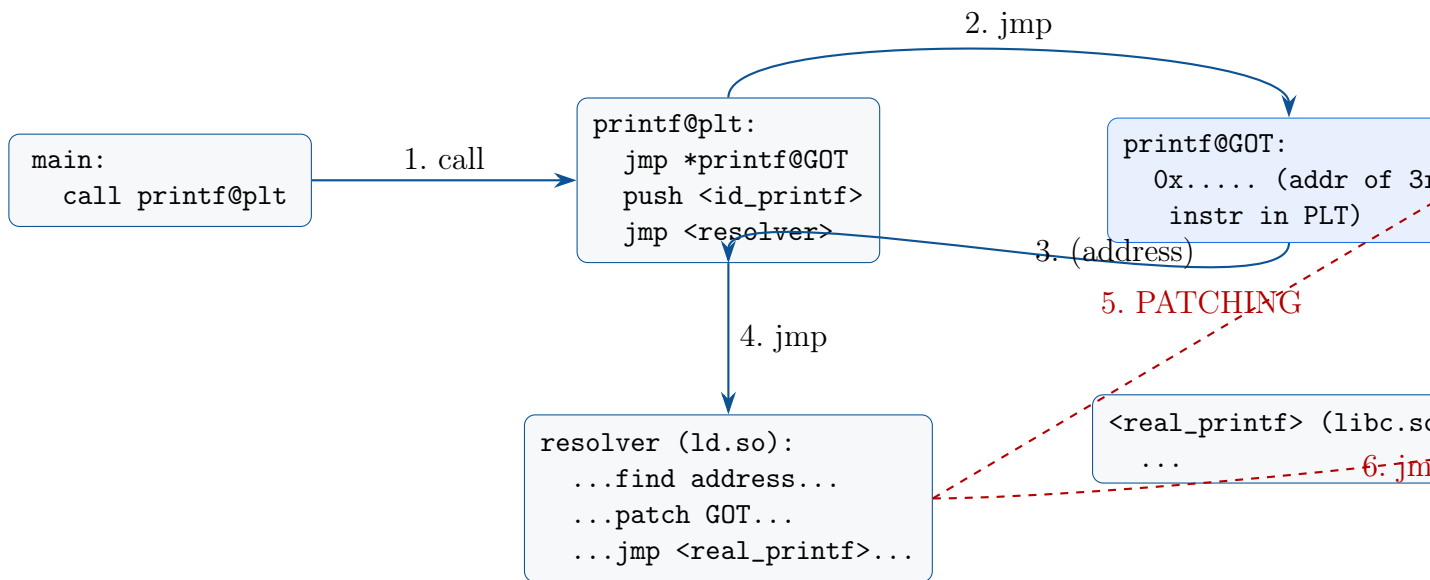


Схема **ленивое связывание** при *первом* вызове

Рис. 8.2 – Процесс ленивого связывания (Lazy Binding)

8.4.3 Перехват вызовов (LD_PRELOAD)

`LD_PRELOAD` — это переменная окружения Linux, которая указывает загрузчику `ld.so`, какую **динамическая библиотека (.so)** загрузить *в первую очередь*, до `libc` и всех остальных.

Если мы создадим свою `libmyhack.so`, в которой определим функцию `printf`, и запустим программу:

```
$ LD_PRELOAD=./libmyhack.so ./a.out
```

Когда резолвер `ld.so` будет искать `printf`, он сначала найдёт *нашу* реализацию в `libmyhack.so` и использует её.

Примечание

Компилятор может заменять небезопасные функции (как `printf`) на их "проверяющие" аналоги (например, `__printf_chk`) для защиты от переполнения буфера. Если вы хотите перехватить `printf`, вам, возможно, придётся перехватывать `__printf_chk`.

8.4.4 Ручная загрузка библиотек (dlopen)

Программа может сама загружать `.so` во время выполнения, используя API из `libdl`.

- `dlopen(const char* path, int mode)`: Загружает `.so`. Возвращает `void*` "handle".
- `dlsym(void* handle, const char* symbol)`: Ищет символ (функцию или переменную) по имени в загруженной библиотеке. Возвращает `void*`.
- `dlclose(void* handle)`: Выгружает библиотеку.
- `dlerror()`: Возвращает строку с описанием последней ошибки.

```
1  #include <dlfcn.h>
2  #include <stdio.h>
3
4  // 1. Define the signature of the function we're looking for
5  typedef double (*sin_func_t)(double);
6
7  int main() {
8      // 2. Load the library
9      void* handle = dlopen("libm.so.6", RTLD_LAZY);
10     if (!handle) { /* error handling */ }
11
12     // 3. Find the symbol (function)
13     void* sym = dlsym(handle, "sin");
14     if (!sym) { /* error handling */ }
15
16     // 4. Cast void* to the correct FUNCTION POINTER type
17     sin_func_t my_sin = (sin_func_t)sym;
18
19     // 5. Use it
20     double result = my_sin(1.0); // ~0.841
21     printf("sin(1.0) = %f\n", result);
22
23     // 6. Close it
24     dlclose(handle);
25     return 0;
26 }
```

Листинг 8.9 – Ручная загрузка `libm.so` для вызова `sin`

Примечание

[title=Внимание!] Ошибка в сигнатуре функции при касте `dlsym` (листинг 8.9, шаг 4) — это тяжёлое **Undefined Behavior**. Если `sin` ожидает `double`, а вы вызовете его с `int`, это почти гарантированно приведёт к падению из-за нарушения соглашения о вызовах (аргументы будут лежать не в тех регистрах, XMM0 vs RDI).

Итоги раздела

- Динамические библиотеки (`.so`) экономят память и позволяют подменять реализации.
- Они должны быть скомпилированы как `PIC` (`-fPIC`) для относительной адресации.
- `PLT` и `GOT` — механизмы, позволяющие вызывать функции, адреса которых неизвестны до запуска.
- `ленивое связывание` (по умолчанию) разрешает адрес функции при первом вызове через `PLT`.
- `LD_PRELOAD` позволяет перехватывать вызовы, подгружая свою `.so` первой.
- `dlopen` и `dlsym` позволяют программе вручную загружать плагины (`.so`) во время работы.

8.5 Freestanding: Программы без `stdlib`

Мы научились делать **Системный вызов** сами. Теперь мы можем полностью отказаться от стандартной библиотеки C/C++.

8.5.1 hosted vs freestanding

- **Hosted:** Стандартный режим. Компилятор предполагает наличие ОС и `stdlib`. Доступны `main`, `malloc`, `printf`, `std::vector` и т.д.
- **Freestanding:** Режим, в котором не предполагается наличие `stdlib`. Нельзя использовать `malloc`, I/O, исключения, RTTI (если они требуют поддержки `stdlib`). Это окружение для ядер ОС, драйверов, микроконтроллеров.

Чтобы собрать программу в `freestanding` режиме, используются флаги:

```
$ g++ -ffreestanding -nostdlib my_program.cpp -o my_program
```

8.5.2 Точка входа `_start`

`main` — это *не* точка входа в программу. Это просто функция, которую вызывает код *инициализации* из `stdlib` (например, `crt0.o`).

Настоящая точка входа, куда ядро Linux передаёт управление, — это метка `_start`. Мы должны определить её сами, обычно на ассемблере.

Состояние при запуске

Когда ядро запускает `_start`, CPU находится в следующем состоянии:

- RSP (указатель стека) 16-байтно выровнен.
- RBP, по соглашению, должен быть обнулён (`xor rbp, rbp`). Это используется дебаггерами и бэктрейсерами как маркер конца цепочки стековых кадров.
- Стек содержит аргументы и переменные окружения (рис. 8.3).

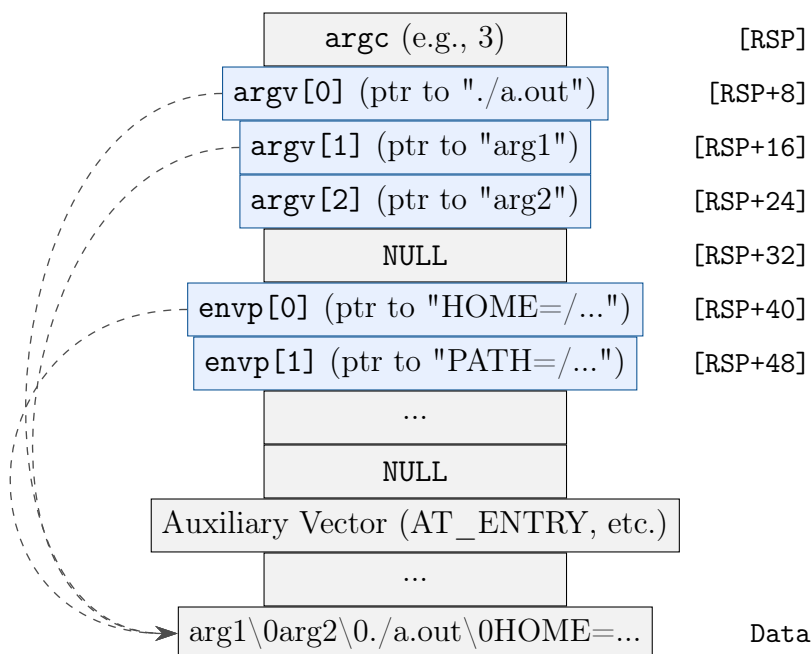


Рис. 8.3 – Содержимое стека при вызове `_start`

Пример `_start`

листинг 8.10 показывает минимальную *freestanding* программу.

```

1  .section .rodata
2  msg:
3      .string "Hello, freestanding!\n"
4  msg_end:
5      .equ msg_len, msg_end - msg
6
7  .section .text
8  .global _start
9
10 _start:
11     # Convention: zero out RBP for backtracing
12     xor %rbp, %rbp
13
14     # syscall: write(1, msg, msg_len)
15     mov $1, %rax # SYS_write
16     mov $1, %rdi # fd (stdout)
17     mov $msg, %rsi # buf
18     mov $msg_len, %rdx # count
19     syscall
20
21     # syscall: exit(123)
22     mov $60, %rax # SYS_exit
23     mov $123, %rdi # exit_code
24     syscall

```

Листинг 8.10 – Freestanding "Hello World" (AT&T синтаксис)

Примечание

Компилятор всё ещё может генерировать вызовы `memcpy`, `memset` и т.д. для оптимизации C++ кода (например, копирования структур). В настоящей *freestanding* среде вам пришлось бы предоставить реализации и этих функций.

8.5.3 Загрузка и расширение знака

При работе с ассемблером важно помнить о размерах данных. Ошибка из прошлой лекции: чтение 32-битного `int` в 64-битный регистр.

- `mov %eax, [%rsp]` (AT&T): Загружает 32 бита из памяти в `EAX`. При этом старшие 32 бита `RAX` обнуляются.
- `mov %rax, [%rsp]`: Загружает 64 бита.

Проблема: если мы читаем 32-битное отрицательное число (`0xFFFFFFFF`, т.е. -1) с помощью `mov %eax, ..., RAX` станет `0x00000000FFFFFFFF` (положительное число ~ 4 млрд).

Определение: Инструкции расширения знака

- `movsx` / `movzx` (Move with Sign Extend) / `movslq` (Move Sign-extend Long to Quad): Копирует знаковый бит (старший бит) источника во все старшие биты приёмника.
- `movzx` (Move with Zero Extend) / `movzb/w...`: Заполняет старшие биты приёмника

нулями.

```
56 # Load a 32-bit dword from [rsp] into 64-bit rax
57 # with sign extension.
58 # Intel: movsx rax, dword ptr [rsp]
59 # AT&T:
60 movslq (%rsp), %rax
```

Итоги раздела

- freestanding режим позволяет писать код без `stdlib`.
- Точка входа в ELF — это `start`, а не `main`.
- Ядро передаёт `argc`, `argv` и `envp` через стек.
- `_start` должна обнулить RBP (`xor %rbp, %rbp`).
- При загрузке 32-битных знаковых чисел в 64-битные регистры необходимо использовать `movsx` / `movzx` (`movslq`) для сохранения знака.

8.6 Введение в архитектуру процессора

Мы научились генерировать инструкции, но почему они выполняются так быстро? Современные CPU — это сложные системы, скрывающие огромную задержку (latency) доступа к памяти.

8.6.1 Проблема доступа к памяти и кэши

1. **Виртуальная память:** Разыменование указателя (виртуального адреса) требует 4-5 обращений к памяти для прохода по таблицам страниц (Page Tables).
2. **DRAM:** Обращение к основной памяти (DRAM) занимает ~ 100 наносекунд.

Итого ~ 500 нс (тысячи тактов CPU) на *каждый* доступ к памяти.

Решение 1: TLB

Translation Lookaside Buffer (буфер ассоциативной трансляции) (TLB) — это маленький, очень быстрый кэш внутри CPU, который хранит недавние отображения "виртуальная страница \rightarrow физическая страница". Если в TLB есть попадание (hit), CPU избегает 4-5 обращений к памяти.

Решение 2: Иерархия кэшей L1/L2/L3

TLB решает проблему трансляции, но кэш решает проблему медленной DRAM.

- **L1 (32-64KB):** ~ 1 нс (несколько тактов). Обычно разделён на L1i (инструкции) и L1d (данные).
- **L2 (256KB-4MB):** $\sim 4-12$ нс.
- **L3 (8MB+):** $\sim 30-50$ нс. (Общий для всех ядер).
- **DRAM:** $\sim 100+$ нс.

Организация кэша

- **кэш-линия:** Данные передаются не побайтно, а блоками по 64 байта.
- **ассоциативный кэш:** Кэш организован как хэш-таблица.

Физический адрес разбивается на три части: [TAG (36 бит) | SET_INDEX (6-10 бит) | OFFSET (6 бит)]

1. **OFFSET (0-5):** Байт внутри 64-байтной кэш-линии.
2. **SET_INDEX (6-11):** Индекс "корзины"(set) в хэш-таблице.
3. **TAG (12-47):** Уникальный идентификатор кэш-линии.

Примечание

Простая хэш-функция (средние биты адреса) — это проблема. Если программа часто обращается к адресам с шагом, кратным большой степени двойки (например, `arr[i * 1024]`), все эти обращения могут попасть в *один и тот же set*, "выбивая" друг друга из кэша, даже если кэш L1 в целом пуст.

8.6.2 Конвейер инструкций (Pipeline)

Для сокращения задержек (даже L1) CPU исполняет инструкции в конвейер (например: Fetch \rightarrow Decode \rightarrow Execute \rightarrow Memory \rightarrow Writeback). Исполнение нескольких инструкций перекрывается во времени.

Конфликты (Hazards)

- **конфликт по данным:** Инструкция N (напр. `add`) ждёт результат инструкции N-1 (напр. `mov`). Пример: итерация по связному списку (`node = node->next`) — это чистый **конфликт по данным**, т.к. следующий `mov` зависит от предыдущего `mov`.
- **конфликт по управлению:** Условный переход (`je`, `jne`). CPU не знает, какую инструкцию загружать (Fetch) следующей, пока не выполнится (Execute) `cmp`.

Продвинутые оптимизации CPU

1. **Out-of-Order Execution (внеочередное исполнение) (OoOE):** CPU может исполнять инструкции не по порядку, если они не зависят друг от друга, чтобы "заполнить" простои (например, во время **конфликт по данным** или промаха кэша).
2. **Переименование регистров:** CPU имеет сотни *физических* регистров, но только 16 *архитектурных* (`rax...`). CPU динамически переименовывает `rax` в `phys_reg_5` в одной инструкции и в `phys_reg_28` в другой, чтобы разорвать ложные зависимости по данным.
3. **предсказание ветвлений:** Для решения **конфликт по управлению**, CPU *угадывает* результат `je` и спекулятивно исполняет код.

Пример: Предсказатель ветвлений

Рассмотрим код (даже на Python) для подсчёта элементов в массиве:

```

1 import numpy as np
2 data = np.random.randint(0, 256, size=1000000)
3 # data.sort() # <--- The key line
4
5 count = 0
6 for x in data:
7     if x < 128: # <--- The conditional branch
8         count += 1

```

Листинг 8.11 – Тест предсказателя ветвлений

- **Несортированный массив:** `if x < 128` непредсказуем. Предсказатель ошибается в ~50% случаев.
- **Сортированный массив:** `if` всегда `True` для первой половины, всегда `False` для второй. Предсказатель ошибается *только один раз* (когда `True` меняется на `False`).

Результат: код на сортированном массиве работает *значительно* (в 5-10 раз) быстрее из-за почти 100% точности **предсказание ветвлений**.

Итоги раздела

- Доступ к DRAM очень медленный (~100 нс).
- TLB кэширует трансляцию виртуальных адресов в физические.
- Кэши L1/L2/L3 кэшируют сами данные из DRAM.
- Данные ходят **кэш-линия** по 64 байта.
- **конвейер** перекрывает исполнение инструкций.

- OoOE, переименование регистров и [предсказание ветвлений](#) — ключевые техники CPU для сокрытия задержек и решения [конфликт по данным](#) и [конфликт по управлению](#).

Глава 9

9 Лекция

9.1 Оптимизации в современных процессорах

Современные CPU применяют множество сложных оптимизаций для достижения высокой производительности. Рассмотрим ключевые из них: организацию кэш-памяти, внеочередное и спекулятивное исполнение инструкций, а также предсказание ветвлений.

9.1.1 Ассоциативность кэша и её влияние на производительность

Кэш — это небольшая, но очень быстрая память, расположенная близко к вычислительным ядрам процессора. Она хранит копии часто используемых данных из основной, более медленной памяти. Эффективность кэша напрямую влияет на скорость работы программ.

Определение: Ассоциативность кэша

Ассоциативность определяет, в скольких возможных местах (слотах) кэша может быть размещена определённая строка данных (кэш-линия) из основной памяти. Кэш-линии группируются в множества (sets). В N -ассоциативном кэше каждая кэш-линия может быть помещена в любое из N мест внутри своего множества.

Типичные значения ассоциативности: 2, 4, 8 или 16. Прямо-отображаемый кэш (1-ассоциативный) прост, но страдает от коллизий: две кэш-линии, претендующие на одно и то же место, будут постоянно вытеснять друг друга. Увеличение ассоциативности снижает вероятность коллизий, но усложняет аппаратуру.

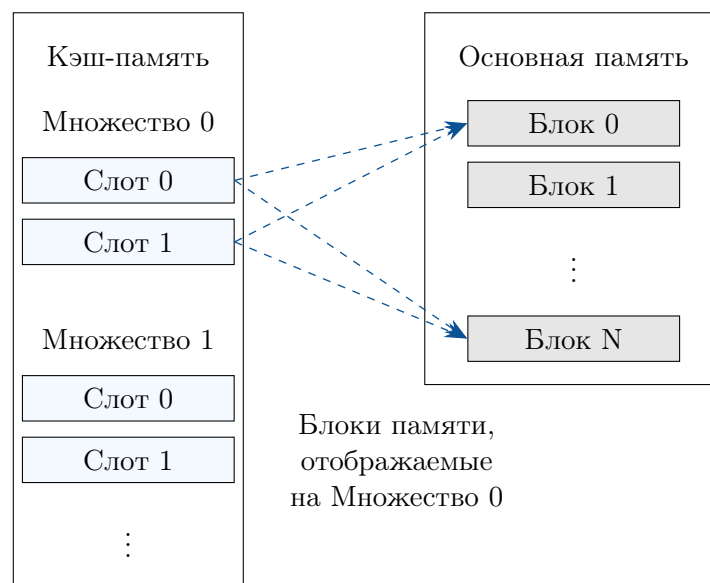


Рис. 9.1 — Схема 2-ассоциативного кэша: любой блок памяти, чей адрес отображается на Множество 0, может быть помещён в любой из двух слотов этого множества.

Неправильный паттерн доступа к памяти может привести к «отравлению» кэша. Рассмотрим пример транспонирования матрицы. При обходе матрицы по столбцам адреса соседних элементов отстоят друг от друга на размер строки. Если размер строки кратен большой степени двойки, адреса элементов из разных строк, но одного столбца, могут отображаться на одно и то же или на малое подмножество множеств в кэше.

Это приводит к постоянным промахам (cache miss), так как кэш-линии вытесняют друг друга. Эксперименты показывают, что транспонирование матрицы 512×512 (где $512 = 2^9$) выполняется значительно медленнее, чем матриц 511×511 или 513×513 , именно по этой причине.

9.1.2 Конвейерное и внеочередное исполнение

Для ускорения обработки инструкций **CPU** использует **конвейер**. Выполнение каждой инструкции разбивается на стадии (выборка, декодирование, исполнение, доступ к памяти, запись результата). Это позволяет одновременно обрабатывать несколько инструкций на разных стадиях.

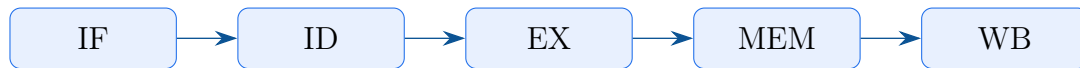


Рис. 9.2 – Классический 5-стадийный конвейер обработки инструкций

Современные процессоры идут дальше и реализуют **OoOE**.

Определение: Внеочередное исполнение (Out-of-Order Execution)

Это способность **CPU** исполнять инструкции не в том порядке, в котором они указаны в программе, а в порядке готовности их операндов. Это позволяет обходить задержки (например, при ожидании данных из памяти) и лучше загружать исполнительные устройства процессора.

Процессор анализирует зависимости по данным между инструкциями. Если две инструкции не зависят друг от друга, они могут быть выполнены параллельно или в обратном порядке. Для разрешения конфликтов по регистрам используется **переименование регистров**: архитектурным регистрам (видимым программисту) ставятся в соответствие физические регистры внутри **CPU**. Это позволяет устранить ложные зависимости.

9.1.3 Спекулятивное исполнение и уязвимости

OoOE тесно связано со спекулятивным исполнением. Процессор может не только переупорядочивать, но и «угадывать» результат условных переходов (ветвлений) и начинать выполнять инструкции из наиболее вероятной ветки кода ещё до того, как условие будет вычислено.

Примечание

Спекулятивное исполнение может оставлять следы в кэше. Если процессор спекулятивно выполнил чтение из памяти, к которой у программы нет доступа, данные могут попасть в кэш. Хотя результат операции будет отброшен после обнаружения ошибки доступа, наличие данных в кэше можно определить по времени доступа к ним. На этом принципе были основаны уязвимости класса **Meltdown** и **Spectre**.

9.1.4 Предсказание ветвлений (Branch Prediction)

Эффективность спекулятивного исполнения зависит от точности предсказания ветвлений. Ошибка предсказания (branch misprediction) очень дорога: **CPU** должен сбросить **конвейер**, отменить результаты спекулятивно выполненных инструкций и начать выполнение с правильной ветки.

Рассмотрим пример: подсчёт элементов в массиве, которые меньше определённого порога.

- **Отсортированный массив:** Предсказатель легко угадывает результат сравнения. Сначала все элементы будут меньше порога, потом — больше. Переход будет только один. Производительность высокая.

- **Неотсортированный (случайный) массив:** Результат сравнения непредсказуем. Процент ошибок предсказания высок ($\approx 50\%$), что приводит к значительному падению производительности.

Компиляторы знают об этой проблеме и могут применять оптимизации, чтобы избежать ветвлений. Например, условное приращение счётчика 'if ($x < 128$) sum++;' может быть заменено на инструкцию условного перемещения (conditional move), которая не содержит прыжка и не нагружает предсказатель ветвлений.

Итоги раздела

- **Ассоциативность кэша** помогает бороться с коллизиями, но паттерны доступа к памяти с шагом, кратным степени двойки, могут снизить её эффективность.
- **Конвейер** и **ОоОЕ** позволяют исполнять несколько инструкций параллельно, скрывая задержки.
- **Спекулятивное исполнение** на основе предсказания ветвлений ускоряет код, но ошибки предсказания дорого обходятся.
- Компиляторы могут преобразовывать код для минимизации ветвлений и улучшения производительности.

9.2 Представление нецелых чисел

Целочисленные типы не могут представлять дробные значения. Для этого в вычислительной технике используются два основных подхода: числа с фиксированной и с плавающей запятой.

9.2.1 Числа с фиксированной точкой (Fixed-Point)

Идея проста: хранить число как целое, но считать, что дробная точка находится в заранее определённой позиции. Фактически это целое число, делённое на фиксированную степень двойки.

- **Преимущества:** Арифметика быстрая, так как используются целочисленные операции.
- **Недостатки:** Ограниченный и фиксированный диапазон значений. Сложно представлять одновременно очень большие и очень маленькие числа. Точность постоянна по всему диапазону.

Например, число 5.125_{10} в двоичном виде равно 101.001_2 . Если мы договоримся хранить 3 знака после запятой, то это число будет храниться как целое 101001_2 .

9.2.2 Стандарт IEEE 754: числа с плавающей запятой

Для гибкого представления широкого диапазона чисел был разработан стандарт [стандарт двоичной арифметики с плавающей запятой \(IEEE 754\)](#). Число представляется в научном формате:

$$fp = S \cdot M \cdot 2^E \quad (9.2.1)$$

где:

- S — знак (+1 или -1).
- M — мантисса (значащая часть), нормализованное число в диапазоне $[1.0, 2.0)$.
- E — экспонента (показатель степени).

В двоичном представлении это выглядит так:

Знак (1 бит)	Экспонента (несколько бит)	Мантисса (остальные биты)
--------------	----------------------------	---------------------------

Поскольку нормализованная мантисса всегда начинается с единицы (1...), эта единица не хранится явно («скрытый бит»), что даёт дополнительный бит точности.

Для хранения отрицательных экспонент используется **смещение (bias)**. Хранимое значение экспоненты — это беззнаковое целое, из которого вычитается bias для получения реального показателя степени.

$$E_{\text{real}} = E_{\text{stored}} - \text{bias} \quad (9.2.2)$$

9.2.3 Специальные случаи

Стандарт [IEEE 754](#) определяет кодирование для особых значений:

- **Денормализованные числа:** Если все биты экспоненты равны 0, скрытый бит считается равным 0 (а не 1). Это позволяет плавно представлять числа, очень близкие к нулю, заполняя «дыру» между нулём и наименьшим нормализованным числом.
- **Бесконечность ($\pm\infty$):** Если все биты экспоненты равны 1, а все биты мантиссы равны 0. Получается при переполнении или делении на ноль ($1.0/0.0$).

- **Не-число** («не число» (**Not a Number**) (**NaN**)): Если все биты экспоненты равны 1, а мантисса не равна нулю. Результат некорректных операций, таких как $\infty - \infty$ или $\sqrt{-1}$.

Примечание

NaN обладает особым свойством: любое сравнение с NaN, даже 'NaN == NaN', возвращает 'false'. Это требует особой осторожности при проверках.

9.2.4 Погрешности и работа в C++

Арифметика с плавающей запятой неточна. Это приводит к нарушению привычных математических законов:

- **Неассоциативность сложения:** $(a + b) + c$ может не равняться $a + (b + c)$, особенно если числа сильно различаются по величине.
- **Недистрибутивность:** $a \cdot (b + c)$ может не равняться $a \cdot b + a \cdot c$.

Для минимизации ошибок при суммировании большого количества чисел их рекомендуется сортировать и складывать от меньших по модулю к большим.

В C++ есть три основных типа с плавающей запятой:

Таблица 9.1 – Типы данных с плавающей запятой в C++

Тип	Размер (байты)	Биты экспоненты	Биты мантиссы
float	4	8	23
double	8	11	52
long double	10	15	64

Для доступа к битовому представлению числа можно использовать 'reinterpret_cast', 'std::bit_cast' (в C++ 20) или структуры с битовыми полями, помня об обратном порядке полей на little-endian архитектурах.

```

1  #include <cstdint>
2
3  // Order is reversed for little-endian systems
4  struct DoubleBits {
5      uint64_t mantissa : 52;
6      uint64_t exponent : 11;
7      uint64_t sign : 1;
8  };
9
10 double d = 1.234;
11 // In C++20, prefer std::bit_cast
12 DoubleBits bits = *reinterpret_cast<DoubleBits*>(&d);
13 // Now bits.sign, bits.exponent, bits.mantissa can be accessed

```

Листинг 9.1 – Доступ к битам double через структуру с битовыми полями

Итоги раздела

- Числа с **фиксированной точкой** просты и быстры, но имеют ограниченный диапазон.
- Стандарт **IEEE 754** определяет представление чисел с **плавающей запятой** (знак, экспонента, мантисса), позволяя работать с огромным диапазоном значений.
- Существуют специальные значения: **денормализованные числа**, **бесконечности** и **NaN**.
- Арифметика с плавающей запятой неточна и требует аккуратного обращения для минимизации погрешностей.

9.3 Основы многопоточности

Многопоточность — это способ организации вычислений, при котором программа состоит из нескольких потоков управления, выполняющихся параллельно.

9.3.1 Процессы и потоки

Определение: Процесс и Поток

Процесс — это экземпляр программы, выполняемый операционной системой. Процессы сильно изолированы друг от друга: у каждого своё адресное пространство, свои файловые дескрипторы и т.д. Коммуникация между ними сложна (требуется IPC: pipes, shared memory).

Поток (thread) — это минимальная единица исполнения внутри процесса. Все потоки одного процесса разделяют общее адресное пространство, файловые дескрипторы и другие ресурсы. Это делает коммуникацию между ними простой, но создаёт проблемы с синхронизацией.

В C++ для создания потоков используется класс `std::thread`.

```
1 #include <iostream>
2 #include <thread>
3
4 void worker_function() {
5     std::cout << "Worker thread is running.\n";
6 }
7
8 int main() {
9     std::thread t(worker_function); // Create and start a new thread
10    // ... main thread continues execution ...
11    t.join(); // Wait for the worker thread to finish
12    return 0;
13 }
```

Листинг 9.2 – Создание и запуск потока в C++

9.3.2 Синхронизация и доступ к общей памяти

Основная сложность в многопоточном программировании — корректная работа с общими данными. Когда несколько потоков одновременно читают и пишут в одну и ту же ячейку памяти, возникает **состояние гонки** (**race condition**).

Проблема усугубляется тем, что и компилятор, и процессор могут переупорядочивать операции для оптимизации. В однопоточной программе это незаметно, но в многопоточной может привести к непредсказуемому поведению.

Примечание

Одновременный доступ (хотя бы одна из операций — запись) к обычной (неатомарной) переменной из разных потоков без синхронизации является **неопределённым поведением** (**неопределённое поведение (Undefined Behavior) (UB)**) в C++.

9.3.3 Атомарные операции (`std::atomic`)

Для безопасной работы с разделяемыми переменными без блокировок используются атомарные типы (`std::atomic`).

Определение: Атомарная операция

Это операция, которая выполняется как единое, неделимое целое. Никакой другой поток не может наблюдать её в промежуточном состоянии.

Например, операция ‘value++’ неатомарна. Она состоит из трёх шагов: чтение, инкремент, запись. Другой поток может вмешаться между этими шагами. Атомарная операция ‘value.fetch_add(1)’ выполняет то же самое, но гарантированно неделимо.

```
1  #include <atomic>
2  #include <thread>
3  #include <vector>
4
5  std::atomic<int> counter = 0;
6
7  void increment() {
8      for (int i = 0; i < 1000000; ++i) {
9          counter.fetch_add(1); // Atomic increment
10     }
11 }
12
13 int main() {
14     std::vector<std::thread> threads;
15     for (int i = 0; i < 10; ++i) {
16         threads.emplace_back(increment);
17     }
18     for (auto& t : threads) {
19         t.join();
20     }
21     // counter will be exactly 10,000,000
22     return 0;
23 }
```

Листинг 9.3 – Безопасный инкремент с помощью std::atomic

9.3.4 Прimitives блокирующей синхронизации

Когда требуется защитить не одну переменную, а целый блок кода (критическую секцию), используются блокирующие примитивы.

Мьютекс (std::mutex)

Мьютекс обеспечивает взаимное исключение. Только один поток может владеть мьютексом в любой момент времени.

- ‘mutex.lock()’: Захватывает **мьютекс**. Если он уже захвачен другим потоком, текущий поток блокируется («засыпает») до его освобождения.
- ‘mutex.unlock()’: Освобождает **мьютекс**.

Для безопасного использования рекомендуется RAII-обёртка ‘std::lock_guard’, которая автоматически вызывает ‘unlock’ в своём деструкторе.

Спинлок (Spinlock)

Альтернатива мьютексу, реализованная на атомарных операциях. Вместо блокировки потока (передачи управления ядру), спинлок входит в цикл активного ожидания (busy-

wait), постоянно проверяя, не освободился ли ресурс.

- **Эффективен**, когда ожидание короткое (меньше, чем накладные расходы на переключение контекста потока).
- **Расточителен**, если ожидание долгое, так как впустую тратит процессорное время.

Условные переменные (std::condition_variable)

Позволяют одному потоку ждать, пока не выполнится некоторое условие, которое устанавливается другим потоком. Они работают в паре с мьютексом.

- 'cv.wait(lock, predicate)': Атомарно освобождает **мьютекс** ('lock') и блокирует поток до тех пор, пока другой поток не вызовет 'notify' и 'predicate' не станет истинным. Перед выходом из 'wait' **мьютекс** снова захватывается.
- 'cv.notify_one()': «Будит» один из ожидающих потоков.

Использование предиката в 'wait' обязательно для борьбы с «ложными пробуждениями» (spurious wakeups).

Итоги раздела

- Потоки разделяют память, что требует **синхронизации** для избежания гонок и **UB**.
- **Атомарные операции** ('std::atomic') обеспечивают неделимый доступ к одиночным переменным.
- **Мьютекс** ('std::mutex') защищает критические секции кода, блокируя потоки при ожидании.
- **Условные переменные** ('std::condition_variable') позволяют потокам эффективно ожидать выполнения произвольных условий.

9.4 Классическая проблема: обедающие философы

Эта задача иллюстрирует проблему **взаимоблокировка** в системах с разделяемыми ресурсами.

9.4.1 Постановка задачи

Пять философов сидят за круглым столом. Перед каждым — тарелка спагетти, а между каждыми двумя соседними философами лежит по одной вилке. Итого 5 философов и 5 вилок.

Каждый философ попеременно то думает, то ест. Чтобы поесть, ему нужны обе вилки: левая и правая.

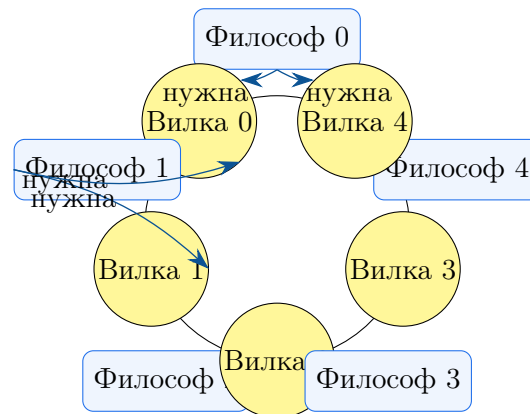


Рис. 9.3 – Схема расположения философов и вилок. Каждому философу для еды нужны две соседние вилки.

9.4.2 Взаимоблокировка (Deadlock)

Рассмотрим наивный алгоритм поведения для каждого философа:

1. Взять левую вилку.
2. Взять правую вилку.
3. Поесть.
4. Положить левую вилку.
5. Положить правую вилку.
6. Подумать.

Примечание

Что произойдёт, если все философы одновременно решат поесть и каждый возьмёт свою левую вилку? Каждый из них будет вечно ждать, пока его сосед справа освободит правую для него вилку. Но сосед справа тоже ждёт. Возникает **цикл ожидания**, и ни один из потоков не может продолжить выполнение. Это и есть **взаимоблокировка**.

Проблема **взаимоблокировка** — одна из фундаментальных в многопоточном программировании. Для её решения существуют различные подходы, например, нарушение одного из условий возникновения взаимоблокировки (в данном случае, введение строгого порядка захвата ресурсов: например, все философы сначала берут вилку с меньшим номером, а потом с большим).

Глава 10

10 Лекция

10.1 Введение: Сигналы и наблюдаемый параллелизм

В начале лекции был затронут вопрос о сигналах: существует ли для них очередь?

Примечание

По умолчанию **очереди сигналов не существует**. Операционная система поддерживает маску ожидающих (pending) сигналов. Если приходит два одинаковых сигнала (например, от пользователя и от системы), а обработчик ещё не вызван, то в маске просто выставляется бит. Второй сигнал того же типа может быть потерян.

10.1.1 Наблюдаемый параллелизм и квантование

Даже на одноядерных системах пользователи наблюдают иллюзию параллельного исполнения множества потоков. Это достигается за счёт **Time Slicing** (квантования времени).

- Операционная система нарезает исполнение потоков на небольшие интервалы — *кванты* (единицы или десятки миллисекунд).
- По истечении кванта таймер прерывает исполнение, планировщик ОС сохраняет контекст текущего потока и загружает следующий.
- Для пользователя переключение происходит незаметно, создавая видимость одновременной работы.

Однако увеличение количества потоков сверх физических возможностей процессора не даёт прироста производительности, а лишь увеличивает накладные расходы на переключение контекста.

10.2 Аппаратная поддержка многопоточности

10.2.1 Hyper-threading

Процессоры часто простаивают в ожидании данных из памяти (cache miss, ~500 тактов) или при выполнении долгих арифметических операций. Чтобы утилизировать эти ресурсы, была разработана технология [Hyper-threading](#).

Определение: Hyper-threading

Технология, позволяющая одному *физическому* ядру процессора представляться операционной системе как два *логических* ядра.

- **Дублируются:** регистры (архитектурное состояние).
- **Разделяются:** исполнительные устройства (ALU), кэши, шина.

Если один логический поток простаивает (ждет память), физическое ядро переключается на инструкции второго потока. Это увеличивает пропускную способность (throughput), но производительность одного отдельного потока может снизиться из-за конкуренции за ресурсы ядра.

10.2.2 Привязка к ядрам (CPU Affinity)

Мы можем программно управлять тем, на каких ядрах выполняется поток, используя системные вызовы семейства `sched_`.

Пример закрепления потока за нулевым ядром:

```
1  cpu_set_t cpuset;  
2  CPU_ZERO(&cpuset);  
3  CPU_SET(0, &cpuset); // Allow execution only on Core 0  
4  
5  pthread_setaffinity_np(thread.native_handle(), sizeof(cpu_set_t), &cpuset);
```

Если запустить два вычислительно тяжёлых потока и привязать их к одному физическому ядру, время выполнения увеличится вдвое по сравнению с их запуском на разных ядрах.

10.2.3 Закон Амдала

Существует теоретический предел ускорения программы при распараллеливании. Он описывается [Закон Амдала](#).

Пусть:

- P — доля программы, которую можно распараллелить (parallelizable).
- S — доля программы, исполняемая последовательно (serial), $S = 1 - P$.
- T — количество потоков.

Ускорение при использовании T потоков:

$$\text{Speedup}(T) = \frac{1}{S + \frac{P}{T}} \quad (10.2.1)$$

Максимально возможное ускорение (при $T \rightarrow \infty$):

$$\text{Max Speedup} = \frac{1}{S} = 1 + \frac{P}{S} \quad (10.2.2)$$

Даже если $P = 0.95$ (95% кода параллелится), максимальное ускорение не превысит 20 раз.

10.3 Thread Safety (Потокобезопасность)

Основная проблема многопоточности — доступ к общим данным.

Определение: Конфликтующая операция (Conflicting Operation)

Операция над одной ячейкой памяти, где участвуют как минимум два потока, и хотя бы один из них выполняет **запись**.

Стандарт C++ гласит: конфликтующие операции, выполняемые одновременно без синхронизации, приводят к **Data Race** и Undefined Behavior.

10.3.1 Контракты стандартной библиотеки C++

Для контейнеров (например, `std::vector`, `std::map`) действуют следующие правила:

1. **Константные методы** ('const') считаются потокобезопасными (можно вызывать одновременно из разных потоков). Например: `'v.size()'`, `'v.capacity()'`.
2. **Неконстантные методы** не являются потокобезопасными. Нельзя вызывать `'v.push_back()'` одновременно с `'v.size()'` или другим `'push_back()'`.
3. Разные экземпляры объектов независимы.

Исключения (thread-safe mutable):

- `std::atomic` — все методы потокобезопасны.
- `std::mutex` — `lock()`/`unlock()` меняют состояние, но безопасны.
- `std::cin`, `std::cout` — операторы `«`, `»` потокобезопасны (не упадут, но вывод может перемешаться).

10.3.2 Свободные функции

- `printf`, `scanf`: Thread-safe (внутри есть глобальные блокировки).
- `getenv`: Помечена как *MT-Safe env*. Безопасна, если никто одновременно не вызывает модифицирующие функции вроде `setenv`, `clearenv`.

10.4 Примитивы синхронизации

Помимо `std::mutex` и `std::condition_variable`, существуют специализированные примитивы.

10.4.1 Семафор (Semaphore)

Счётчик разрешений (permits). Не имеет владельца (в отличие от мьютекса, который должен освобождать тот же поток, что и захватил).

- **Wait (acquire):** Декрементирует счётчик. Если 0 — блокирует поток.
- **Signal (release):** Инкрементирует счётчик, будит ждущий поток.

Применение: ограничение количества одновременных доступов (rate limiting), очереди. В C++20: `std::counting_semaphore`, `std::binary_semaphore`.

10.4.2 RWLock (Read-Write Lock)

Позволяет множеству читателей работать одновременно, но писателю даёт эксклюзивный доступ. В C++: `std::shared_mutex`.

- `lock_shared()`: Захват для чтения (несколько потоков могут держать одновременно).
- `lock()`: Захват для записи (эксклюзивно: ни читателей, ни писателей).

Проблема: Возможен *starvation* (голодание) писателей, если поток читателей непрерывен, или читателей, если приоритет у писателей.

10.4.3 Барьер (Barrier)

Синхронизирует прохождение N потоками определённой точки. Пример: нейросети. Нужно посчитать слой i полностью всеми потоками, прежде чем переходить к слою $i + 1$.

```
1 // Example of using std::barrier
2 std::barrier sync_point(num_threads);
3
4 // In the worker thread code:
5 for (int l = 0; l < layers; ++l) {
6     ProcessLayerPart(l);
7     sync_point.arrive_and_wait(); // Wait for others
8 }
```

10.5 Thread Local Storage (TLS)

Иногда каждому потоку нужна своя копия глобальной переменной (например, буфер или код ошибки). Используется ключевое слово `thread_local`.

```
thread_local int value; // Each thread has its own instance
```

Классический пример: `errno`. В однопоточных системах это была глобальная переменная. В многопоточных она реализована как макрос, возвращающий разыменованный указатель на TLS-переменную:

```
#define errno (*__errno_location())
```

Это позволяет разным потокам иметь разные коды ошибок одновременно.

Итоги раздела

Итоги по инструментам:

- Используйте `const` методы для параллельного чтения.
- Для специфичных задач используйте `semaphore` (лимиты), `shared_mutex` (read-heavy), `barrier` (этапы).
- `thread_local` для изоляции данных потока.

10.6 Shared Pointer и безопасность

Вопрос: является ли `std::shared_ptr` потокобезопасным? **Ответ:** Частично.

- **Control Block (счётчик ссылок):** Потокобезопасен. Инкремент/декремент счётчика атомарен (используются атомарные инструкции). Можно копировать/удалять 'shared_ptr' из разных потоков.
- **Указываемый объект:** НЕ защищён. Если два потока пишут в данные через разные 'shared_ptr', указывающие на один объект — будет гонка.

Для отладки гонок полезен инструмент **ThreadSanitizer (TSan)**. Он ловит:

- Write-Read races.
- Write-Write races.

Запуск: компиляция с '-fsanitize=thread'.

10.7 Атомики и модель памяти

10.7.1 Compare-And-Swap (CAS)

Мощная атомарная операция для реализации lock-free алгоритмов. В C++: `compare_exchange_strong` и `compare_exchange_weak`.

Логика работы:

```
1 // CAS Pseudocode
2 bool CAS(atomic<int>& obj, int& expected, int desired) {
3     if (obj == expected) {
4         obj = desired;
5         return true;
6     } else {
7         expected = obj; // Updates 'expected' with actual value
```

```

8     return false;
9   }
10 }

```

Использование в цикле (CAS-loop) для реализации 'fetch_add':

```

1 std::atomic<int> atom;
2 int expected = atom.load();
3 while (!atom.compare_exchange_weak(expected, expected + 1)) {
4     // If failed, 'expected' is automatically updated inside the function
5     // Loop repeats with new 'expected'
6 }

```

Weak vs Strong:

- **Weak:** Может вернуть 'false' (спонтанный сбой), даже если значение равно ожидаемому. На платформах вроде ARM/PowerPC это эффективнее (нет вложенного цикла). На x86 (Intel/AMD) разницы в ассемблере нет (обе транслируются в 'LOCK CMPXCHG').
- **Strong:** Гарантирует успех, если значение совпало. Обычно содержит цикл внутри себя на RISC-архитектурах.

10.8 Кэши и когерентность

10.8.1 Иерархия памяти и протокол MESI

У каждого ядра есть свои L1/L2 кэши. Если одно ядро пишет в переменную, другие должны узнать об этом, чтобы не читать устаревшие данные. Для этого используется протокол когерентности, например, **MESI**:

- **M (Modified):** Данные изменены, есть только в этом кэше.
- **E (Exclusive):** Данные только в этом кэше, совпадают с памятью.
- **S (Shared):** Данные есть в нескольких кэшах (только чтение).
- **I (Invalid):** Данные устарели.

Коммуникация между ядрами для поддержки когерентности стоит дорого.

10.8.2 False Sharing (Ложное разделение)

Проблема возникает, когда независимые переменные попадают в одну **кэш-линию** (обычно 64 байта).

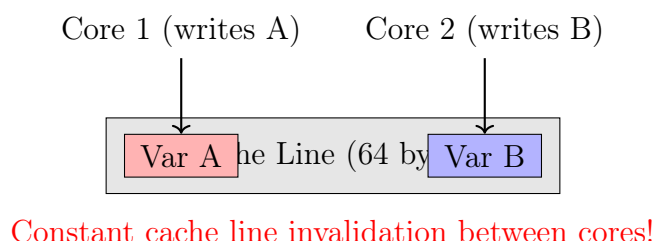


Рис. 10.1 – Иллюстрация False Sharing

Если Thread 1 пишет в 'Var A', а Thread 2 пишет в 'Var B', и они лежат рядом, ядра будут бесконечно передавать друг другу права на владение всей кэш-линией. Это катастрофически снижает производительность.

Решение: Выравнивание (padding).

```
1 struct alignas(64) AlignedData {
2     int value;
3     // Compiler adds padding up to 64 bytes
4 };
```

10.9 Процессы и Fork в многопоточной среде

Системный вызов `fork()` создаёт копию процесса. **Правило:** В дочернем процессе продолжает исполнение *только тот поток*, который вызвал `fork`. Остальные потоки исчезают.

Проблема: Если исчезнувший поток держал мьютекс (например, внутри `'malloc'`), этот мьютекс останется навечно заблокированным в дочернем процессе. Любая попытка вызвать `'malloc'` в "ребёнке" приведёт к дедлоку.

Решение: Использовать `'pthread_atfork'` для регистрации хуков, которые захватывают нужные блокировки перед форком и отпускают их после (в обоих процессах), обеспечивая консистентное состояние.

10.10 Futex (Fast Userspace Mutex)

Системный вызов Linux, на котором строятся эффективные мьютексы. Идея: избегать входа в ядро (syscall), если нет конкуренции (fast path). В ядро идем только чтобы уснуть (wait).

- `futex_wait(addr, val)`: "Если по адресу `addr` лежит значение `val`, то усыпи меня". Проверка атомарна внутри ядра.
- `futex_wake(addr, count)`: Разбуди `count` потоков, ждущих на этом адресе.

Пример примитивного мьютекса на атомике и фьютексе:

```
1 std::atomic<int> flag{0}; // 0 - free, 1 - locked
2
3 void lock() {
4     while (flag.exchange(1) != 0) { // Try to acquire
5         // If busy, go to wait state
6         syscall(SYS_futex, &flag, FUTEX_WAIT, 1, ...);
7     }
8 }
9
10 void unlock() {
11     flag.store(0);
12     syscall(SYS_futex, &flag, FUTEX_WAKE, 1, ...);
13 }
```

Примечание

В реальных реализациях (`'std::mutex'`) используется более сложная логика с тремя состояниями (свободен, занят, занят+есть ждущие), чтобы избежать лишних системных вызовов `'wake'`.

Итоги раздела**Итоги лекции:**

1. Параллелизм ограничивается не только числом ядер, но и синхронизацией (Закон Амдала) и аппаратными эффектами (False Sharing).
2. Hyper-threading позволяет утилизировать простаивающие ядра, но делит ресурсы.
3. Вызов `fork()` в многопоточной программе опасен из-за состояния блокировок.
4. Современные блокировки строятся на атомике + Futex для эффективного ожидания.

Глава 11

11 Лекция

Оглавление

11.1 Введение

В рамках данной лекции рассматривается стек сетевых технологий: от физических принципов передачи сигнала до реализации высоконагруженных серверов с использованием механизмов мультиплексирования ввода-вывода. Основное внимание уделяется архитектурным ограничениям (trade-offs) различных протоколов и системных вызовов.

11.2 Физический и Канальный уровни

11.2.1 Принципы передачи сигнала

Фундаментальной задачей сетевого взаимодействия является передача информации между двумя физически соединёнными узлами. На физическом уровне это реализуется посредством модуляции напряжения в проводнике. Выделяют два типа сигналов:

1. **Аналоговый сигнал.** Обладает непрерывным спектром значений. Позволяет передавать большой объём информации, однако критически подвержен зашумлению, что делает его непригодным для точной передачи данных в современных компьютерных сетях.
2. **Цифровой сигнал.** Использует дискретный набор значений напряжения (например, высокий и низкий уровень). Значения, выходящие за пределы порогов, округляются, что обеспечивает устойчивость к помехам.

При передаче последовательности бит (например, длинной серии единиц) возникает проблема синхронизации: получатель (Receiver) может рассинхронизироваться с отправителем (Sender) относительно частоты тактов. Для решения данной проблемы применяется **Манчестерское кодирование**.

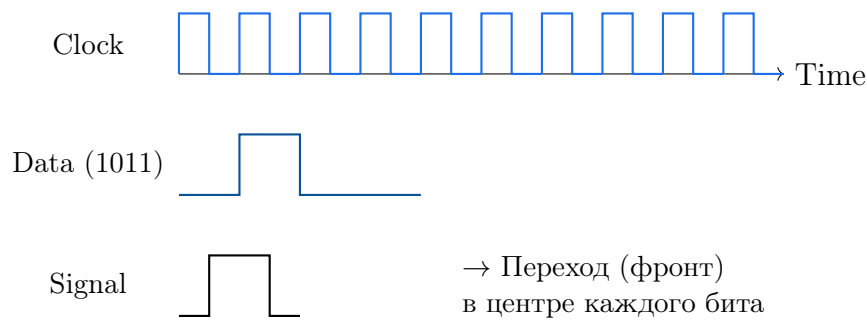
Определение: Манчестерское кодирование

Метод кодирования, при котором данные логически объединяются (операция XOR) с тактовым сигналом (Clock). Это гарантирует наличие перепада напряжения (фронта) в середине каждого битового интервала, что позволяет получателю синхронизировать частоту по самому сигналу.

11.2.2 Протокол Ethernet и MAC-адресация

Для структурирования потока бит используется протокол [Ethernet](#), оперирующий единицами данных, называемыми фреймами (frames).

Адресация на канальном уровне осуществляется посредством [MAC-адрес](#)-адресов. Предполагается глобальная уникальность MAC-адресов, обеспечиваемая распределением диапазонов (OUI) между производителями оборудования (например, Intel, Cisco).

**Рис. 11.1** – Принцип Манчестерского кодирования (схематично)**Таблица 11.1** – Структура Ethernet-фрейма

Поле	Описание	Размер
Preamble	Синхронизация	8 байт
Destination MAC	MAC-адрес получателя	6 байт
Source MAC	MAC-адрес отправителя	6 байт
EtherType	Тип инкапсулированного протокола (IPv4/IPv6)	2 байта
Payload	Полезная нагрузка (например, IP-пакет)	46-1500 байт
FCS	Контрольная сумма (CRC) для проверки целостности	4 байта

11.3 Сетевой уровень (Network Layer)

11.3.1 Протокол IP и маршрутизация

Протокол **IP** (Internet Protocol) обеспечивает глобальную адресацию и маршрутизацию пакетов между сетями. Наиболее распространённая версия IPv4 использует 32-битные адреса, записываемые в виде четырёх октетов (например, 192.168.0.1).

Ввиду ограниченности адресного пространства IPv4 (всего $2^{32} \approx 4.3$ млрд адресов), широко применяется технология **NAT**.

Определение: NAT (Network Address Translation)

Механизм, позволяющий устройствам из локальной сети (с приватными адресами) выходить в глобальную сеть, используя один публичный IP-адрес маршрутизатора. Маршрутизатор подменяет адрес источника и порт, сохраняя состояние соединения в таблице трансляции.

11.3.2 Подсети и маски (CIDR)

Для логического разделения сетей используется бесклассовая адресация (CIDR). Адрес сети определяется префиксом, а маска подсети указывает количество бит, отведённых под адрес сети.

```

1 IP: 88.99.146.0
2 CIDR: /23 (23 bits network, 9 bits host)
3
4 Netmask: 11111111.11111111.11111110.00000000 (255.255.254.0)
5 Wildcard: 00000000.00000000.00000001.11111111 (0.0.1.255)
6
7 Network: 88.99.146.0
8 HostMin: 88.99.146.1

```

```

9 HostMax: 88.99.147.254
10 Broadcast: 88.99.147.255
11 Hosts/Net: 510 (2^9 - 2)

```

Листинг 11.1 – Пример расчета диапазона подсети /23

Примечание

Адреса, у которых хостовая часть состоит из всех нулей (адрес сети) или всех единиц (Broadcast), зарезервированы и не могут быть назначены конкретному интерфейсу.

Для предотвращения бесконечной циркуляции пакетов при ошибках маршрутизации используется поле **TTL** (Time To Live), значение которого уменьшается на единицу при прохождении каждого промежуточного узла (хопа).

11.4 Транспортный уровень (Transport Layer)

Сетевой уровень доставляет данные до хоста. Транспортный уровень (L4) отвечает за мультиплексирование потоков данных для конкретных приложений, используя абстракцию **порта** (16-битное число).

11.4.1 Сравнение TCP и UDP

Таблица 11.2 – Сравнение транспортных протоколов

Характеристика	TCP (Stream)	UDP (Datagram)
Гарантии	Гарантирует доставку и порядок байт.	Не гарантирует доставку и порядок.
Соединение	Требуется установка соединения (handshake).	Отправка данных без установки соединения ("fire and forget").
Поток данных	Непрерывный поток байт.	Дискретные сообщения (датаграммы).
Overhead	Высокий (заголовки, подтверждения).	Низкий (минимальный заголовок).
Применение	Web (HTTP), Email, File transfer.	DNS, Streaming, VoIP, Gaming.

11.4.2 Порядок байт (Endianness)

Сетевые протоколы исторически используют порядок байт **Big-Endian** (старший байт по младшему адресу). Архитектура x86 использует **Little-Endian**. Для корректной интерпретации многобайтовых чисел (например, порта или IP-адреса) необходима конвертация:

```

1 #include <arpa/inet.h>
2
3 uint16_t host_port = 8080;
4 // Host TO Network Short (conversion for sending)
5 uint16_t net_port = htons(host_port);
6 // Network TO Host Short (conversion on receipt)
7 uint16_t back_port = ntohs(net_port);

```

Листинг 11.2 – Конвертация порядка байт

11.5 Прикладной уровень (Application Layer)

Прикладные протоколы (L7) определяют семантику передаваемых данных.

11.5.1 HTTP и DNS

Протокол HTTP — текстовый протокол, работающий поверх TCP. Сообщения состоят из стартовой строки, заголовков и опционального тела. В качестве разделителя строк используется последовательность CRLF (\r\n).

Система DNS выполняет трансляцию доменных имён (например, google.com) в IP-адреса. Для получения адреса в языке C используется функция `getaddrinfo`, возвращающая список структур `addrinfo`.

Примечание

Один домен может разрешаться (resolve) в несколько IP-адресов для балансировки нагрузки на уровне DNS (DNS Round Robin) и обеспечения географической близости к клиенту (CDN).

11.6 Socket API и Модели конкурентности

11.6.1 Базовый жизненный цикл TCP-сервера

Интерфейс сокетов (сокет) в POSIX-системах реализует файловую абстракцию для сетевого взаимодействия.

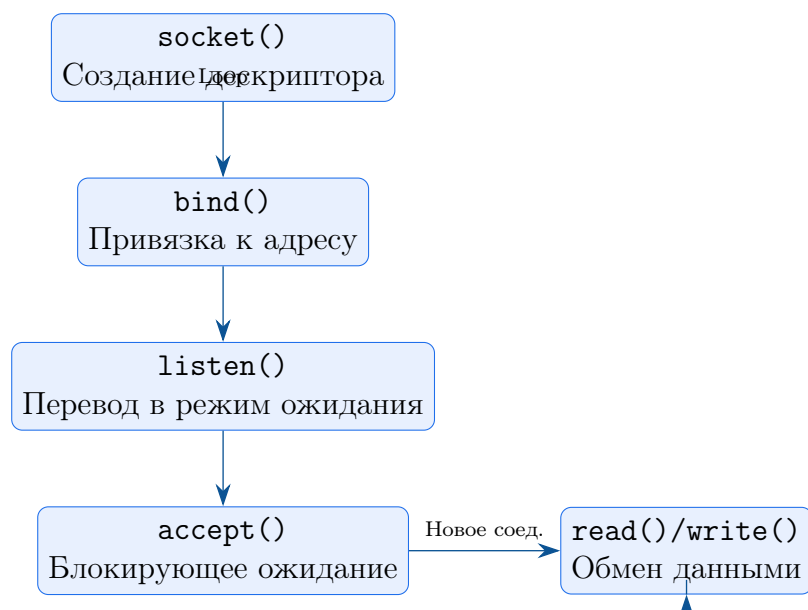


Рис. 11.2 – Системные вызовы TCP-сервера

Критическая проблема базовой модели: вызов `read()` является блокирующим. Если сервер обслуживает одного клиента, остальные клиенты ожидают в очереди (backlog), создаваемой ядром ОС.

11.6.2 Эволюция моделей ввода-вывода

Thread-per-connection (Поток на соединение)

Создание отдельного потока ОС (pthread) для каждого клиента.

- **Преимущество:** Простота реализации, линейный код.

- **Цена (Trade-off):** Высокие накладные расходы. Каждый поток требует аллокации стека (по умолчанию 2-8 МБ) и структур ядра. Переключение контекста (Context Switch) при тысячах потоков приводит к существенной деградации производительности (CPU тратится на планировщик, а не на полезную работу).

Non-blocking I/O (Неблокирующий ввод-вывод)

Файловый дескриптор переводится в неблокирующий режим (`O_NONBLOCK`). Системные вызовы `read/write` возвращают ошибку `EAGAIN`, если данные не готовы. Это позволяет одному потоку опрашивать множество сокетов, но приводит к проблеме **Busy Wait** (100% загрузка CPU при пустых циклах опроса).

I/O Multiplexing (epoll)

Механизм `epoll` (Linux) позволяет потоку «заснуть» до возникновения события на одном из множества наблюдаемых дескрипторов.

Определение: Event Loop

Архитектурный паттерн, в котором единый поток циклически ожидает событий от мультиплексора (`epoll`) и вызывает соответствующие обработчики (`callbacks`). Это позволяет обрабатывать десятки тысяч соединений (проблема C10k) с минимальными накладными расходами памяти и отсутствием лишних переключений контекста.

```
1 int epfd = epoll_create1(0);
2 struct epoll_event ev, events[MAX_EVENTS];
3
4 // Register server socket
5 ev.events = EPOLLIN;
6 ev.data.fd = server_socket;
7 epoll_ctl(epfd, EPOLL_CTL_ADD, server_socket, &ev);
8
9 while (1) {
10     // Wait for events (blocking)
11     int nfds = epoll_wait(epfd, events, MAX_EVENTS, -1);
12
13     for (int i = 0; i < nfds; ++i) {
14         if (events[i].data.fd == server_socket) {
15             // handle new connection (accept)
16         } else {
17             // handle data (read/write)
18         }
19     }
20 }
```

Листинг 11.3 – Скелет Event Loop на `epoll`

11.7 Итоги раздела

Итоги раздела

- Сетевое взаимодействие представляет собой иерархию абстракций: от модуляции напряжения до прикладных протоколов (OSI Model).

- Протоколы TCP и UDP предоставляют различные гарантии (надёжность против скорости), что диктует их области применения.
- Прямое использование потоков (Thread-per-connection) неэффективно для I/O-bound задач с большим числом соединений из-за накладных расходов на память и планирование.
- Современные высоконагруженные серверы используют неблокирующий ввод-вывод и механизм [epoll](#) для мультиплексирования событий в едином цикле (Event Loop).

