

Содержание

| | |
|---|-----------|
| Предисловие | 17 |
| Извинения за код | 19 |
| Использование примеров кода | 19 |
| Соглашения, использованные в книге | 20 |
| Об авторе | 20 |
| Об изображении на обложке | 20 |
| Ждем ваших отзывов! | 21 |
| Глава 1. Обзор оптимизации | 23 |
| Оптимизация — часть разработки программного обеспечения | 24 |
| Эффективность оптимизации | 25 |
| Оптимизируйте! | 25 |
| Наносекунда туда, наносекунда сюда | 28 |
| Стратегии оптимизации кода на C++ | 28 |
| Используйте компилятор получше; используйте компилятор лучше | 29 |
| Использование лучших алгоритмов | 30 |
| Использование лучших библиотек | 32 |
| Уменьшение количества выделений памяти и копирований | 33 |
| Устранение вычислений | 33 |
| Использование лучших структур данных | 34 |
| Увеличение параллельности | 34 |
| Оптимизация управления памятью | 35 |
| Резюме | 35 |
| Глава 2. Оптимизация, влияющая на поведение компьютера | 37 |
| Ложь о компьютерах, в которую верит C++ | 38 |
| Правда о компьютерах | 39 |
| Медленная память | 40 |
| Недоступность байтов | 41 |
| Одни обращения к памяти медленнее других | 41 |
| Остроконечные и тупоконечные слова | 42 |
| Количество памяти ограничено | 43 |
| Медленное выполнение команд | 43 |
| Трудное принятие решений | 44 |
| Множественные потоки выполнения | 45 |
| Вызовы операционной системы являются дорогостоящими | 46 |
| C++ тоже лжет | 46 |
| Не все инструкции одинаково дорогие | 47 |
| Инструкции выполняются не по порядку | 47 |
| Резюме | 48 |

| | |
|---|-----------|
| Глава 3. Измерение производительности | 49 |
| Оптимизирующее мышление | 50 |
| Производительность должна быть измерена | 50 |
| Оптимизация — большая игра | 51 |
| Правило 90/10 | 51 |
| Закон Амдала | 53 |
| Проведение экспериментов | 54 |
| Ведение лабораторного журнала | 56 |
| Измерение базовой производительности и постановка целей | 57 |
| Улучшить можно только измеряемое | 60 |
| Профилирование выполнения программы | 60 |
| Длительно работающий код | 63 |
| “Полузнайство” об измерении времени | 64 |
| Измерение времени с помощью компьютеров | 69 |
| Преодоление проблем измерений | 77 |
| Создание класса-секундомера | 81 |
| Хронометраж функции в тесте | 85 |
| Оценка стоимости кода для поиска узких мест | 86 |
| Оценка стоимости отдельных инструкций C++ | 87 |
| Оценка стоимости циклов | 88 |
| Другие пути поиска узких мест | 89 |
| Резюме | 90 |
| Глава 4. Оптимизация использования строк | 91 |
| Почему строки представляют собой проблему | 91 |
| Строки используют динамическое выделение памяти | 92 |
| Строки как значения | 92 |
| Строки выполняют массу копирований | 93 |
| Первая попытка оптимизации строк | 94 |
| Использование модифицирующих операций для устранения временных значений | 96 |
| Уменьшение работы с памятью с помощью резервирования | 96 |
| Устранение копирования строкового аргумента | 97 |
| Устранение разыменований с помощью итераторов | 98 |
| Устранение копирования возвращаемого значения | 99 |
| Использование массивов символов вместо строк | 100 |
| Итоги первой попытки оптимизации | 102 |
| Вторая попытка оптимизации строк | 102 |
| Использование лучшего алгоритма | 102 |
| Использование лучшего компилятора | 104 |
| Использование лучшей библиотеки для работы со строками | 105 |
| Использование лучшего менеджера памяти | 109 |

| | |
|---|------------|
| Устранение преобразования строк | 110 |
| Преобразование строк в стиле C в <code>std::string</code> | 111 |
| Преобразование между кодировками | 111 |
| Резюме | 112 |
| Глава 5. Оптимизация алгоритмов | 113 |
| Временная стоимость алгоритмов | 115 |
| Временная стоимость в наилучшем, среднем и наихудшем случаях | 117 |
| Амортизированная временная стоимость | 118 |
| Прочие стоимости | 118 |
| Оптимизации сортировки и поиска | 118 |
| Эффективные алгоритмы поиска | 119 |
| Временная стоимость алгоритмов поиска | 119 |
| Все поиски равноценны при малых n | 120 |
| Эффективные алгоритмы сортировки | 121 |
| Временная стоимость алгоритмов сортировки | 121 |
| Замена сортировки с плохой производительностью в наихудшем случае | 122 |
| Использование информации о входных данных | 123 |
| Шаблоны оптимизации | 123 |
| Предвычисления | 124 |
| Отложенные вычисления | 125 |
| Пакетирование | 126 |
| Кеширование | 126 |
| Специализация | 127 |
| Группировка | 127 |
| Подсказки | 128 |
| Оптимизация ожидаемого пути | 128 |
| Хеширование | 128 |
| Двойная проверка | 129 |
| Резюме | 129 |
| Глава 6. Оптимизация переменных в динамической памяти | 131 |
| Переменные C++ | 132 |
| Длительность хранения переменной | 132 |
| Владение переменными | 135 |
| Объекты-значения и объекты-сущности | 136 |
| API динамических переменных C++ | 138 |
| Автоматизация владения интеллектуальными указателями | 140 |
| Динамические переменные имеют стоимость времени выполнения | 143 |
| Уменьшение использования динамических переменных | 144 |
| Статическое создание экземпляров класса | 145 |
| Использование статических структур данных | 146 |
| Использование <code>std::make_shared</code> вместо <code>new</code> | 150 |

| | |
|---|------------|
| Не следует разделять владение без необходимости | 150 |
| Использование “главного указателя” для владения динамическими переменными | 152 |
| Уменьшение количества перераспределений динамических переменных | 152 |
| Предварительное выделение памяти для динамических переменных для предотвращения перераспределений | 152 |
| Создание динамических переменных вне циклов | 153 |
| Устранение излишнего копирования | 154 |
| Устранение нежелательного копирования в определении класса | 155 |
| Устранение копирования при вызове функции | 156 |
| Устранение копирования при возврате из функции | 158 |
| Библиотеки без копирования | 160 |
| Реализация идиомы “копирования при записи” | 161 |
| Срезы | 162 |
| Реализация семантики перемещения | 163 |
| Нестандартная семантика копирования: болезненный хак | 163 |
| std::swap(): семантика перемещения для бедных | 164 |
| Разделяемое владение сущностями | 165 |
| Перемещающая часть семантики перемещения | 166 |
| Изменения кода для использования семантики перемещения | 167 |
| Тонкости семантики перемещения | 168 |
| Плоские структуры данных | 171 |
| Резюме | 172 |
| Глава 7. Оптимизация инструкций | 173 |
| Удаление кода из циклов | 174 |
| Кеширование конечного значения цикла | 175 |
| Применение более эффективных инструкций циклов | 175 |
| Изменение направления цикла | 176 |
| Устранение инвариантного кода из циклов | 177 |
| Удаление ненужных вызовов функций из циклов | 177 |
| Удаление скрытых вызовов функций из циклов | 180 |
| Удаление дорогих медленно меняющихся вызовов из циклов | 182 |
| Перемещение циклов в функции для снижения накладных расходов при вызовах | 183 |
| Выполняйте некоторые действия пореже | 184 |
| И все остальное | 186 |
| Удаление кода из функций | 186 |
| Стоимость вызовов функций | 186 |
| Объявление коротких функций встраиваемыми | 190 |
| Определение функций до их первого использования | 191 |
| Устранение неиспользуемого полиморфизма | 191 |
| Удаление неиспользуемых интерфейсов | 192 |
| Выбор реализации во время компиляции с помощью шаблонов | 196 |
| Исключение применения идиомы PIMPL | 196 |

| | |
|---|------------|
| Устранение вызовов кода в DLL | 198 |
| Используйте статические функции-члены вместо функций-членов экземпляров | 199 |
| Перенесение виртуального деструктора в базовый класс | 199 |
| Оптимизация выражений | 200 |
| Упрощение выражений | 201 |
| Группирование констант | 202 |
| Используйте менее дорогостоящие операторы | 203 |
| Использование целочисленной арифметики вместо арифметики с плавающей точкой | 203 |
| double может быть быстрее, чем float | 205 |
| Замена итеративных вычислений аналитическими выражениями | 206 |
| Идиомы оптимизации потока управления | 207 |
| Применение switch вместо if-elseif-else | 208 |
| Применение виртуальных функций вместо switch или if | 208 |
| Используйте обработку исключений без стоимости | 209 |
| Резюме | 211 |
| Глава 8. Использование лучших библиотек | 213 |
| Оптимизация использования стандартной библиотеки | 213 |
| Философия стандартной библиотеки C++ | 214 |
| Вопросы применения стандартной библиотеки C++ | 215 |
| Оптимизация существующих библиотек | 217 |
| Изменения должны быть небольшими | 218 |
| Добавление функций, а не изменение функциональности | 219 |
| Проектирование оптимизированных библиотек | 219 |
| Кодировать на скорую руку — обречь себя на долгую муку | 219 |
| При разработке библиотек скупость является добродетелью | 221 |
| Принятие решений о выделении памяти вне библиотеки | 221 |
| Если сомневаетесь, выбирайте скорость | 222 |
| Оптимизация функций проще оптимизации каркасов | 222 |
| Плоские иерархии наследования | 223 |
| Упрощение цепочки вызовов | 223 |
| Упрощение проектирования слоев | 223 |
| Избегайте динамического поиска | 225 |
| Остерегайтесь “функций Бога” | 226 |
| Резюме | 227 |
| Глава 9. Оптимизация сортировки и поиска | 229 |
| Таблицы “ключ/значение” с использованием std::map и std::string | 230 |
| Инструментарий для повышения производительности поиска | 231 |
| Выполнение базовых измерений | 232 |
| Идентификация оптимизируемой деятельности | 232 |
| Разделение оптимизируемой деятельности | 233 |

| | |
|---|-----|
| Изменение или замена алгоритмов и структур данных | 234 |
| Использование процесса оптимизации для пользовательских абстракций | 236 |
| Оптимизация поиска с использованием <code>std::map</code> | 237 |
| Применение символьных массивов фиксированного размера в качестве ключей <code>std::map</code> | 237 |
| Использование строк в стиле C в качестве ключей <code>std::map</code> | 238 |
| Использование <code>std::set</code> , когда ключ является значением | 241 |
| Оптимизация поиска с использованием заголовочного файла <code><algorithm></code> | 241 |
| Таблица “ключ/значение” для поиска в последовательных контейнерах | 243 |
| <code>std::find()</code> : очевидное имя, стоимость — $O(n)$ | 244 |
| <code>std::binary_search()</code> : не возвращает значения | 245 |
| Бинарный поиск с использованием <code>std::equal_range()</code> | 245 |
| Бинарный поиск с использованием <code>std::lower_bound()</code> | 246 |
| Самостоятельное кодирование бинарного поиска | 247 |
| Самостоятельное кодирование бинарного поиска с использованием <code>strcmp()</code> | 248 |
| Оптимизация поиска в хешированных таблицах “ключ/значение” | 248 |
| Хеширование с использованием <code>std::unordered_map</code> | 249 |
| Хеширование с фиксированными символьными массивами в качестве ключей | 250 |
| Хеширование с ключами в виде строк с завершающими нулевыми символами | 251 |
| Хеширование с пользовательской хеш-таблицей | 253 |
| Цена абстракций Степанова | 254 |
| Оптимизация сортировки с использованием стандартной библиотеки C++ | 255 |
| Резюме | 257 |

Глава 10. Оптимизация структур данных **259**

| | |
|---|-----|
| Знакомство с контейнерами стандартной библиотеки | 259 |
| Последовательные контейнеры | 260 |
| Ассоциативные контейнеры | 260 |
| Эксперименты с контейнерами стандартной библиотеки | 261 |
| <code>std::vector</code> и <code>std::string</code> | 266 |
| Следствия перераспределения для производительности | 267 |
| Вставка и удаление в <code>std::vector</code> | 268 |
| Итерирование <code>std::vector</code> | 270 |
| Сортировка <code>std::vector</code> | 271 |
| Поиск в <code>std::vector</code> | 271 |
| <code>std::deque</code> | 271 |
| Вставка и удаление в <code>std::deque</code> | 273 |
| Итерирование <code>std::deque</code> | 275 |
| Сортировка <code>std::deque</code> | 275 |
| Поиск в <code>std::deque</code> | 275 |
| <code>std::list</code> | 275 |
| Вставка и удаление в <code>std::list</code> | 278 |
| Итерирование <code>std::list</code> | 278 |

| | |
|--|------------|
| Сортировка <code>std::list</code> | 278 |
| Поиск в <code>std::list</code> | 279 |
| <code>std::forward_list</code> | 279 |
| Вставка и удаление в <code>std::forward_list</code> | 280 |
| Итерирование <code>std::forward_list</code> | 280 |
| Сортировка <code>std::forward_list</code> | 281 |
| Поиск в <code>std::forward_list</code> | 281 |
| <code>std::map</code> and <code>std::multimap</code> | 281 |
| Вставка и удаление в <code>std::map</code> | 282 |
| Итерирование <code>std::map</code> | 284 |
| Сортировка <code>std::map</code> | 285 |
| Поиск в <code>std::map</code> | 285 |
| <code>std::set</code> и <code>std::multiset</code> | 285 |
| <code>std::unordered_map</code> и <code>std::unordered_multimap</code> | 286 |
| Вставка и удаление в <code>std::unordered_map</code> | 289 |
| Итерирование <code>std::unordered_map</code> | 290 |
| Поиск в <code>std::unordered_map</code> | 290 |
| Другие структуры данных | 290 |
| Резюме | 292 |
| Глава 11. Оптимизация ввода-вывода | 293 |
| Рецепты для чтения файлов | 293 |
| Создание экономной сигнатуры функции | 295 |
| Сокращение цепочек вызовов | 297 |
| Снижение количества перераспределений | 297 |
| Использование большего входного буфера | 299 |
| Использование построчного чтения | 300 |
| Еще одно сокращение цепочек вызовов | 302 |
| Бесполезные вещи | 303 |
| Запись файлов | 303 |
| Чтение из <code>std::cin</code> и запись в <code>std::cout</code> | 304 |
| Резюме | 305 |
| Глава 12. Оптимизация параллельности | 307 |
| Введение в параллельные вычисления | 308 |
| Экскурсия по зоопарку параллелизма | 309 |
| Чередующееся выполнение | 313 |
| Последовательная согласованность | 314 |
| Гонки | 315 |
| Синхронизация | 316 |
| Атомарность | 317 |
| Возможности параллельности в C++ | 319 |
| Потоки | 320 |
| Обещания и фьючерсы | 321 |

| | |
|---|------------|
| Асинхронные задания | 323 |
| Мьютексы | 325 |
| Блокировки | 326 |
| Условные переменные | 327 |
| Атомарные операции над общими переменными | 330 |
| Будущие возможности параллелизма C++ | 333 |
| Оптимизация многопоточных программ C++ | 334 |
| Предпочитайте <code>std::async</code> , а не <code>std::thread</code> | 335 |
| Создавайте потоков столько же, сколько имеется ядер | 337 |
| Реализуйте очередь заданий и пул потоков | 338 |
| Выполняйте ввод-вывод в отдельном потоке | 339 |
| Программа без синхронизации | 339 |
| Удаление кода запуска и завершения | 342 |
| Более эффективная синхронизация | 343 |
| Уменьшайте критические разделы | 344 |
| Ограничивайте количество параллельных потоков | 344 |
| Избегайте громового стада | 346 |
| Избегайте очереди на блокировку | 346 |
| Уменьшайте конкуренцию | 347 |
| Не пользуйтесь активным ожиданием в одноядерных системах | 348 |
| Не ждите вечно | 349 |
| Собственные мьютексы могут быть неэффективными | 349 |
| Ограничивайте длину очереди вывода производителя | 349 |
| Библиотеки для параллельных вычислений | 350 |
| Резюме | 351 |
| Глава 13. Оптимизация управления памятью | 353 |
| API управления памятью C++ | 354 |
| Жизненный цикл динамических переменных | 354 |
| Функции для выделения и освобождения памяти | 355 |
| Построение динамических переменных с помощью выражений <code>new</code> | 358 |
| Уничтожение динамических переменных с помощью выражения <code>delete</code> | 361 |
| Явный вызов деструктора уничтожает динамическую переменную | 362 |
| Высокопроизводительные диспетчеры памяти | 363 |
| Диспетчеры памяти для конкретных классов | 365 |
| Диспетчер памяти для блоков фиксированного размера | 366 |
| Арена блоков | 369 |
| Добавление <code>operator new()</code> для конкретного класса | 371 |
| Производительность диспетчера памяти для блоков фиксированного размера | 372 |
| Вариации диспетчера памяти для блоков фиксированного размера | 372 |
| Небезопасные с точки зрения параллельности диспетчеры более эффективны | 373 |

| | |
|--|------------|
| Пользовательские аллокаторы стандартной библиотеки | 374 |
| Минимальный аллокатор в C++11 | 376 |
| Дополнительные определения для аллокатора C++98 | 378 |
| Аллокатор блоков фиксированного размера | 382 |
| Аллокатор блоков фиксированного размера для строк | 384 |
| Резюме | 385 |
| Предметный указатель | 387 |

Оптимизация, влияющая на поведение компьютера

Ложь, умение рассказывать прекрасные истории, каких никогда не случилось, составляет истинную цель Искусства.¹

— Оскар Уайльд (Oscar Wilde), “Упадок лжи”, *Намерения* (1891)

Цель настоящей главы — обеспечить минимум справочной информации об оборудовании компьютера для мотивации оптимизаций, описанных в этой книге, чтобы читателю не пришлось сходить с ума, углубляясь в 600-страничный справочник по процессору. Здесь приводится поверхностный обзор архитектуры процессора, позволяющий выделить некоторые эвристики для оптимизации. Очень нетерпеливый читатель может пропустить эту главу и вернуться к ней позже, когда другие главы будут ссылаться на материал из данной главы. Тем не менее приводимая здесь информация является важной и полезной для понимания остального материала книги.

Микропроцессорные устройства в настоящее время невероятно разнообразны. Они варьируются от дешевых встроенных устройств, состоящих всего лишь из нескольких тысяч логических вентилях с тактовой частотой ниже 1 МГц, до настольных устройств с миллиардами вентилях и частотами, измеряемыми в гигагерцах. Мощные ЭВМ могут иметь размер большой комнаты, содержать тысячи независимых исполнительных устройств и потреблять электроэнергию, достаточную для освещения небольшого города. Заманчиво считать, что ничто не связывает это изобилие вычислительных устройств, но в действительности они имеют очень много общего. В конце концов, если бы между ними не было никакого сходства, было бы невозможно компилировать код на C++ для множества процессоров, для которых имеются соответствующие компиляторы.

Все широко используемые компьютеры выполняют команды, хранящиеся в памяти. Эти команды обрабатывают данные, которые также хранятся в памяти. Память делится на множество небольших *слов* по нескольку битов. Несколько драгоценных слов памяти представляют собой *регистры*, которые именуются в машинных

¹ Перевод А. Зверева.

командах непосредственно. Большинство же слов памяти именуются с использованием их числового *адреса*. Определенный регистр в каждом компьютере содержит адрес следующей выполняемой команды. Если память рассматривать как книгу, то *выполняемый адрес* выглядит как палец, указывающий на следующее читаемое слово. *Исполнительное устройство* (именуемое также процессором, ядром, ЦПУ и кучей других слов) считывает поток команд из памяти и действует в соответствии с ними. Команды говорят исполнительному устройству, какие данные читать (загружать, выбирать) из памяти, что с ними делать и где в памяти записывать (запоминать, сохранять) результаты. Компьютер состоит из устройств, которые подчиняются физическим законам. Чтение или запись каждого адреса памяти занимает некоторое ненулевое количество времени, как и выполнение некоторых действий над считанными данными.

За пределами этой базовой схемы, знакомой любому первокурснику, генеалогическое дерево компьютерных архитектур демонстрирует буйный рост и изобилие ветвей. Поскольку компьютерные архитектуры весьма изменчивы, трудно сформулировать какие-либо строгие числовые правила в отношении поведения аппаратного обеспечения. Современные процессоры выполняют так много различных взаимодействующих операций, чтобы ускорить выполнение команд, что какие-то конкретные сроки их выполнения в общем случае указать невозможно. С учетом того, что многие разработчики даже не знают точно, на каких процессорах будет работать их код, лучшее, что можно ожидать в смысле оптимизации, — это некоторые эвристики.

Ложь о компьютерах, в которую верит C++

Конечно, программа на C++ по крайней мере делает вид, что верит в изложенную в предыдущем разделе версию простой модели компьютера. Существует память, адресуемая в байтах размера `char`, которая является по существу бесконечной. Существует специальный адрес, именуемый `nullptr`, который отличается от любого допустимого адреса памяти. Целое число `0` преобразуется в `nullptr`, хотя `nullptr` не обязательно указывает на адрес `0`. Существует единый концептуальный адрес выполнения, указывающий на инструкцию исходного кода, выполняемую в настоящее время. Инструкции выполняются в том порядке, в котором они написаны, с учетом действия операторов управления потоком выполнения C++.

C++ знает, что в действительности компьютеры сложнее, чем эта простая модель. Поэтому из сияющих внутренностей реализации C++ сквозь крышку пробиваются следующие сверкающие лучи.

- Программа на C++ должна лишь вести себя так, “как если бы” инструкции выполнялись в указанном порядке. Компилятор C++ и сам компьютер имеют право изменять порядок выполнения, чтобы добиться ускорения работы программы, если при этом не меняется смысл выполняемых вычислений.
- Что касается стандарта C++11, то C++ больше не считает, что существует только один адрес выполнения. Стандартная библиотека C++ теперь поставляется с возможностями запуска и остановки потоков выполнения и синхронизации

доступа к памяти между этими потоками. До C++11 программисты лгали компилятору C++ о потоках, что зачастую приводило к трудно отлаживаемым проблемам.

- Некоторые адреса памяти на самом деле вместо обычной памяти могут быть регистрами устройства. Значения в некоторых адресах могут изменяться между двумя последовательными чтениями одного и того же адреса в одном и том же потоке, отражая некоторые действия со стороны аппаратного обеспечения. Такие места описаны в C++ с помощью ключевого слова `volatile`. Объявление переменной как `volatile` требует от компилятора извлекать новую копию переменной всякий раз, когда она используется, вместо того, чтобы оптимизировать программу путем сохранения значения в регистре и его повторного использования. Могут быть также объявлены указатели на `volatile`-память.
- C++11 предлагает магическое заклинание `std::atomic<>`, которое заставляет память некоторое время вести себя так, как если бы это в действительности было простое линейное хранилище байтов, без учета всех сложностей современных микропроцессоров с их многочисленными потоками выполнения, многослойными кешами памяти и т.д. Некоторые разработчики считают, что для этого служит ключевое слово `volatile`, но они сильно ошибаются.

Операционная система также врет программам и их пользователям. На самом деле вся цель операционной системы — сообщить каждой программе набор очень убедительной лжи. Среди наиболее важной лжи — операционная система хочет убедить каждую программу, что она для компьютера единственная, что предоставляемая ей физическая память бесконечна и что есть бесконечное число процессоров, доступных для запуска потоков программы.

Операционная система умело использует аппаратное обеспечение компьютеров, чтобы скрывать свою ложь, так что программа на C++ не имеет иного выбора, кроме как поверить ей. Эта ложь обычно не слишком сильно влияет на работу программы, за исключением замедления ее работы. Однако она может осложнить выполнение измерений производительности программы.

Правда о компьютерах

Модели C++ соответствуют только простейшие микропроцессоры и некоторые древние ЭВМ. Что действительно важно для оптимизации — это то, что аппаратное обеспечение памяти реальных компьютеров очень медленное по сравнению со скоростью выполнения команд, что обращение к памяти в действительности не байтное, что память не является простым линейным массивом одинаковых ячеек и что она имеет конечную емкость. Реальные компьютеры могут иметь более чем один адрес команды. Реальные компьютеры быстрые не потому, что они быстро выполняют каждую команду, а потому, что выполнение нескольких команд перекрывается, и сложные схемы следят за тем, чтобы такие перекрывающиеся команды вели себя так, как если бы они выполнялись одна за другой.

Медленная память

Основная память компьютера очень медленная по сравнению с его внутренними логическими вентилями и регистрами. Она медленная настолько, что процессор настольного компьютера может выполнить сотни команд за время, необходимое для извлечения из основной памяти единственного слова данных.

Следствием этого для оптимизации является то, что *доступ к памяти доминирует над другими действиями процессора*, включая выполнение команд.

Узкое место фон Неймана

Интерфейс к основной памяти является узким местом, которое ограничивает скорость работы. Это узкое место даже имеет имя. Оно называется узким местом (или бутылочным горлышком) фон Неймана, в честь знаменитого пионера компьютерной архитектуры и математика Джона фон Неймана (John von Neumann) (1903–1957).

Например, компьютер, оснащенный памятью DDR2 с частотой 1000 МГц (типичной для компьютеров несколько лет назад), имеет теоретическую пропускную способность 2 млрд слов в секунду, или 500 пс на слово. Но это не значит, что компьютер может считывать или записывать случайное слово данных каждые 500 пс.

Начнем с того, что за один такт (половина такта часов с частотой 1000 МГц) может выполняться только последовательный доступ. Доступ к не последовательным местам в памяти выполняется примерно за 6–10 тактов.

Далее, за доступ к шине памяти конкурируют несколько действий. Процессор постоянно проводит выборку очередных выполняемых команд из памяти. Контроллер кеш-памяти выполняет выборку блоков данных памяти для кеша и сбрасывает записанные строки кеша. Контроллер DRAM также требует циклы для обновления заряда в динамических ячейках устройства памяти. Числа ядер многоядерного процессора достаточно, чтобы гарантировать перенасыщенность шины памяти. Фактическая скорость, с которой данные могут быть прочитаны из основной памяти в определенное ядро, составляет более 20–80 нс на одно слово.

Закон Мура позволяет каждый год получать все больше и больше ядер в процессорах. Но чтобы обеспечить более быстрый интерфейс основной памяти, этого мало. Таким образом, удвоение числа ядер в будущем будет оказывать на производительность отрицательное влияние. Ядра будут бороться за доступ к памяти. Надвигающееся ограничение производительности называется *стеной памяти*.

Недоступность байтов

Хотя C++ считает, что каждый байт доступен по отдельности, компьютеры часто компенсируют медленность физической памяти путем выборки данных большими блоками. Самые мелкие процессоры могут выбирать из основной памяти отдельные байты, но процессоры настольных компьютеров могут выбирать за один раз по 64 байта. Некоторые суперкомпьютеры и графические процессоры выполняют выборку еще большего размера.

Когда C++ выбирает многобайтные данные, такие как тип `int`, `double` или указатель, может оказаться, что байты, составляющие данные, входят в два слова физической памяти. Это называется *невыровненным доступом к памяти*. Важным для оптимизации является то, что *невыровненный доступ требует в два раза больше времени*, чем если бы все байты были в одном и том же слове, поскольку при этом требуется чтение двух слов. Компилятор C++ выравнивает структуры таким образом, чтобы каждое поле начиналось с адреса байта, кратного размеру поля. Но это создает собственную проблему: “дыры” в структурах, содержащие неиспользуемые данные. Обращая внимание на размер полей данных и их порядок в структуре, можно сделать структуры максимально компактными при сохранении выровненности.

Одни обращения к памяти медленнее других

Для компенсации “медленности” основной памяти многие компьютеры содержат *кеш-память*, разновидность быстрой временной памяти, располагающейся близко к процессору, чтобы ускорить доступ к наиболее часто используемым словам памяти. Одни компьютеры обходятся без кеша; другие имеют один или несколько уровней кеша, каждый из которых меньше, быстрее и дороже предыдущего. Когда процессору нужны байты из слова кешированной памяти, они могут быть быстро извлечены без повторного обращения к основной памяти. Насколько кеш-память быстрее? Эмпирическое правило гласит, что каждый уровень кеш-памяти примерно в 10 раз быстрее, чем предыдущий уровень в иерархии памяти. На процессорах настольных компьютеров время доступа к памяти может меняться на пять порядков в зависимости от того, к чему осуществляется доступ: к кеш-памяти первого, второго или третьего уровня, к основной памяти или к странице виртуальной памяти на диске. Это одна из причин, по которым увлеченность тактами выполнения команд процессором и другими подобными тайнами так часто оказывается глупой и бесполезной — состояние кеша делает время выполнения команды весьма неопределенным.

Когда процессор нуждается в выборке данных, находящихся не в кеше, другие данные, расположенные в кеше, могут быть удалены из него для того, чтобы освободить место. Данные, выбираемые для удаления, как правило, являются наиболее давно использовавшимися. Это важно для оптимизации, поскольку означает, что *доступ к интенсивно используемым ячейкам памяти можно получить быстрее, чем к используемым реже*.

Чтение даже одного байта данных, который не находится в кеше, приводит к кешированию множества близлежащих байтов (как следствие это означает, что множество байтов, находящихся в настоящее время в кеше, удаляются из него).

Эти близлежащие байты будут готовы для быстрого доступа. Это важно для оптимизации, поскольку означает, что обращение к соседним ячейкам памяти (в среднем) быстрее, чем к памяти в отдаленных местах.

В терминах C++ это означает, что блок кода, содержащий цикл, может выполняться быстрее, поскольку инструкции, составляющие цикл, активно используются, располагаются близко одна к другой и, таким образом, вероятно, будут оставаться в кеше. Блок кода, содержащий вызовы функций или инструкции `if`, которые приводят к дальним переходам, могут выполняться более медленно, потому что используются части кода, далеко отстоящие одна от другой. Такой код использует больше пространства кеша, чем цикл. Если программа большая, а кеш конечен, часть кода будет удаляться из кеша, чтобы освободить место для других вещей, что приведет к замедлению доступа в следующий раз, когда потребуется этот код. Аналогично доступ к структуре данных, состоящей из последовательных местоположений в памяти, такой как массив или вектор, может быть быстрее, чем к структуре данных, состоящих из узлов, связанных указателями, поскольку данные в последовательных местоположениях с большей вероятностью будут оставаться в кеш-памяти. Доступ к структуре данных, состоящих из записей, связанных с помощью указателей (например, к списку или дереву), может быть медленнее из-за необходимости считывания данных каждого узла из основной памяти в новые строки кеша.

Остроконечные и тупоконечные слова

Из памяти может быть выбран один байт данных, но зачастую одновременно извлекается несколько последовательных байтов, образующих число. Например, в Microsoft Visual C++ значение типа `int` образуют четыре байта, считываемые из памяти вместе. Так как к памяти можно обращаться двумя способами, люди, которые проектируют компьютеры, должны ответить на важный вопрос: что содержится в первом байте (адрес которого наименьший) — старшие или младшие биты значения `int`?

На первый взгляд кажется, что это не должно иметь значения. Конечно, важно, чтобы все части компьютера одинаково считали, с какого конца `int` адрес меньше, иначе воцарится хаос. Вот в чем заключается разница между этими способами хранения. Если значение `int`, равное `0x01234567`, хранится по адресам 1000–1003 и первым хранятся старшие биты, то по адресу 1000 содержится байт `0x01`, а по адресу 1003 — байт `0x67`, в то время как если сначала хранится младший байт, то по адресу 1000 содержится `0x67`, а по адресу 1003 — `0x01`. Компьютеры, которые хранят старшие биты в байте с младшим адресом, называются компьютерами с *обратным порядком байтов* (“тупоконечниками”, *big-endian*). Компьютеры с *прямым порядком байтов* (“остроконечники”, *little-endian*) сначала читают младшие биты. Итак, имеется два способа хранения целого числа (или указателя), и нет никаких причин предпочесть один другому, так что разные команды, работающие на разных процессорах для разных компаний, могут делать разный выбор.

Проблемы начинаются, когда данные, записанные на диск или отправленные по сети одним компьютером, должны быть прочитаны другим компьютером. Диски и сети пересылают информацию побайтно, а не весь `int` одновременно. Поэтому

оказывается важно, какой конец числа сохраняется (или отправляется) первым. Если отправляющий и принимающий компьютеры не согласованы, то значение, отправляемое как 0x01234567, может быть получено как 0x67452301.

Порядок байтов является лишь одной из причин, по которым C++ не определяет, как биты располагаются в `int` или как значение одного поля в объединении влияет на другие поля. Это одна из причин, по которым программу можно написать так, что она будет успешно работать на компьютере одного вида, но приводить к аварийному завершению на другом.

Количество памяти ограничено

Память компьютера не бесконечна. Чтобы сохранить иллюзию бесконечной памяти, операционная система может использовать физическую память наподобие кеш-памяти и хранить данные, которые не помещаются в физическую память, в виде файла на диске. Эта схема называется *виртуальной памятью*. Виртуальная память создает иллюзию большего количества физической памяти. Однако получение блока памяти с диска занимает десятки миллисекунд — вечность для современного компьютера.

Быстрая кеш-память весьма дорогостоящая. В настольном компьютере или смартфоне может быть гигабайт памяти, но кеш размером лишь в несколько мегабайтов. Программы и их данные обычно в кеш не помещаются.

Результатом кеширования и применения виртуальной памяти может быть то, что *из-за кеширования определенная функция, работающая в контексте всей программы, может выполняться медленнее, чем та же функция в тестовой программе*, в 10 тысяч раз. В контексте всей программы функции и данные не могут все время оставаться в кеше, в то время как в контексте теста это именно так и происходит. Этот эффект усиливает преимущества оптимизаций, которые снижают использование памяти или диска, в то время как преимущества оптимизаций, которые уменьшают размер кода, остаются небольшими.

Вторым результатом кеширования является то, что если большая программа выполняет рассеянный доступ ко многим участкам памяти, то кеш-памяти может оказаться недостаточно для хранения данных, непосредственно используемых программой. Это приводит к снижению производительности, которое называют *пробуксовкой страницы*. Когда пробуксовка страницы происходит во внутреннем кеше микропроцессора, результатом является снижение производительности. Когда это происходит в файле виртуальной памяти операционной системы, производительность падает тысячекратно. Эта проблема возникала чаще, когда физическая память была дороже и меньше по размеру, но она встречается и сейчас.

Медленное выполнение команд

Простые микропроцессоры наподобие встроенных в кофеварки и микроволновые печи предназначены для выполнения команд с той же скоростью, с которой они могут извлекаться из памяти. Микропроцессоры настольного компьютера имеют дополнительные ресурсы для параллельной обработки нескольких команд, поэтому они способны выполнять команды во много раз быстрее, чем те могут быть извлечены из

основной памяти, так что большую часть времени для хранения команд и их передачи процессору используется быстрая кеш-память. Важность этого для оптимизации заключается в том, что *время доступа к памяти превышает время вычислений*.

Современные настольные компьютеры выполняют команды с удивительной скоростью, *если* им ничто не мешает. Они могут завершать команды каждые несколько сотен пикосекунд (пикосекунда представляет собой 10^{-12} с, до смешного короткое время). Но это не значит, что каждая команда выполняется так быстро. Процессор содержит “конвейер” одновременно выполняемых команд. Команды проходят через конвейер, дешифруются, получают свои аргументы, выполняют вычисления и сохраняют результаты. Чем более мощный процессор, тем более сложен его конвейер, разбивающий выполнение команды на десятков этапов так, чтобы как можно больше команд могли быть обработаны одновременно.

Если команда А вычисляет значение, которое требуется команде Б, то команда Б не может выполнить свое вычисление до тех пор, пока команда А не даст необходимый результат. Это приводит к *остановке конвейера*, короткой паузе в выполнении команд, которая возникает, когда выполнение двух команд не может полностью перекрываться. Остановка конвейера в особенности долгая, если команда А извлекает значение из памяти, а затем выполняет вычисление, которое дает значение, необходимое команде Б. Остановке конвейера подвержены все современные микропроцессоры, что делает их время от времени почти такими же медленными, как процессор в вашем тостере.

Трудное принятие решений

Еще одна вещь, которая может вызвать остановку конвейера, — это принятие решения компьютером. После большинства команд выполнение продолжается с команды, находящейся в памяти по следующему адресу. Большую часть времени эта следующая команда уже находится в кеше. Последовательные команды можно загружать в конвейер, как только на первом этапе конвейера для этого появляется место.

Но не таковы команды передачи управления. Команда перехода или вызова подпрограммы заменяет адрес выполнения произвольным новым значением. “Следующая” команда не может быть считанной из памяти и попасть в конвейер до тех пор, пока в некоторый момент в процессе обработки команды перехода не будет обновлен адрес выполнения. У слова памяти по новому адресу выполнения меньше шансов находиться в кеше. Конвейер останавливается на время обновления адреса выполнения и загрузки в конвейер новой “следующей” команды.

После команды условного ветвления выполнение продолжается в одном из двух разных мест: выполняется либо следующая команда, либо команда по адресу, который является целевым для команды ветвления, в зависимости от результатов некоторых предыдущих вычислений. Конвейер останавливается, пока не будут завершены все команды, участвующие в предыдущих вычислениях, и остается в этом состоянии до тех пор, пока не будет определен следующий выполняемый адрес и не будет прочитана соответствующая команда по этому адресу.

Важность этого явления для оптимизации заключается в том, что *вычисление быстрее принятия решения*.

Множественные потоки выполнения

Любая программа, работающая в современной операционной системе, разделяет компьютер с другими программами, выполняемыми в то же время, с периодическими обслуживающими процессами типа проверки диска или поиска обновлений Java или Flash, и с различными частями операционной системы, управляющими сетевым интерфейсом, дисками, звуковыми устройствами и другими периферийными устройствами. *Каждая программа конкурирует с другими программами за ресурсы компьютера.*

Программа обычно не слишком об этом осведомлена. Она просто работает не-много медленнее. Исключением является ситуация, когда одновременно запускается много программ, и все они конкурируют за память и диск. Для настройки производительности, *если программа должна запускаться при запуске системы или в периоды пиковой нагрузки, ее производительность должна измеряться под нагрузкой.*

По состоянию на начало 2016 года настольные компьютеры имеют до 16 ядер процессора, а микропроцессоры, используемые в телефонах и планшетах, — до восьми. Беглый взгляд на диспетчер задач Windows, вывод состояния процесса Linux или список задач Android обычно показывает гораздо больше процессов, чем ядер, и большинство процессов имеют несколько потоков выполнения. Операционная система выполняет каждый поток в течение короткого времени, а затем переключает контекст на другой поток или процесс. С точки зрения программы это выглядит, как если бы одна команда выполнялась наносекунду, а следующая — 60 мс.

Что означает переключение контекстов? Если операционная система переключается от одного потока к другому в одной программе, это означает сохранение регистров процессора для приостановки потока и загрузки сохраненных регистров для возобновления другого потока. Регистры современного процессора содержат сотни байтов данных. Когда новый поток возобновляет выполнение, данные могут не быть в кеше, так что имеется начальный период медленного выполнения, пока новый контекст загружается в кеш. Таким образом, при переключении контекстов потоков имеются значительные затраты.

Процедура переключения операционной системой контекста от одной программы к другой еще более дорогостоящая. Все “грязные” страницы кеша (с записанными данными, которые еще не внесены в основную память) должны быть сброшены в физическую память, а все регистры процессора сохранены. Затем сохраняются регистры страниц отображения физической памяти на виртуальную в диспетчере памяти. Далее для нового процесса загружаются соответствующие регистры памяти и регистры процессора. И наконец выполнение программы может возобновиться. Но кеш в этот момент пуст, так что начальный период характеризуется низкой производительностью и конфликтами памяти.

Когда программа должна ожидать некоторое событие, это ожидание может продолжаться даже после того, как это событие произойдет, пока операционная система не освободит процессор для продолжения программы. При выполнении программы в контексте работы других программ, конкурирующих за ресурсы компьютера, это может увеличить время выполнения программы и сделать его более неопределенным.

Исполнительные устройства многоядерных процессоров и связанная с ними кеш-память для достижения лучшей производительности работают более или менее независимо друг от друга. Однако все исполнительные устройства имеют одну и ту же основную память. Они вынуждены конкурировать за доступ к оборудованию, связывая его с основной памятью, что делает узкое место фон Неймана в компьютере с несколькими исполнительными устройствами еще более ограничивающим.

Когда исполнительное устройство записывает значение, оно сначала попадает в кеш-память. В конечном итоге оно должно быть записано из кеша в основную память, так что это значение станет видимым для других исполнительных устройств. Однако из-за конфликтов доступа к основной памяти среди исполнительных устройств основная память может не обновляться в течение сотен команд после того, как значение было изменено.

Если на компьютере имеется несколько исполнительных устройств, то одно из них может, таким образом, на протяжении длительного периода времени не увидеть данные, записанные другим устройством в основной памяти, а изменения в основной памяти могут произойти не в том же порядке, что и порядок выполнения команд. В зависимости от непредсказуемых временных факторов исполнительное устройство может увидеть как старое значение общей памяти, так и обновленное значение. Для того чтобы различные исполнительные устройства видели согласованные представления памяти, должны использоваться специальные команды синхронизации. Значение этого эффекта для оптимизации состоит в том, что *обращение к совместно используемым потоками выполнения данным гораздо медленнее, чем к не разделяемым данным.*

Вызовы операционной системы являются дорогостоящими

Все процессоры, кроме самых мелких, имеют аппаратное обеспечение для обеспечения изоляции между программами, так что программа А не может читать, писать или выполнять команды в физической памяти, принадлежащей программе Б. То же самое оборудование защищает ядро операционной системы от перезаписи программами. С другой стороны, ядру операционной системы требуется доступ к памяти, принадлежащей каждой программе, чтобы эти программы могли делать системные вызовы операционной системы. Некоторые операционные системы также позволяют программам делать запросы для совместного использования памяти. Способы организации системных вызовов и совместно используемой памяти разнообразны и полны тайной магии. С точки зрения оптимизации важным является то, что *системные вызовы дороги*, в сотни раз дороже, чем вызовы функций внутри одного потока одной программы.

C++ тоже лжет

Самая большая ложь C++ в том, что он рассказывает своим пользователям, что компьютер, который выполняет программу, представляет собой простую последовательную структуру. В обмен на то, что разработчики делают вид, что верят в эту ложь, C++ позволяет разработчикам программировать без знаний интимных подробностей каждого микропроцессорного устройства, которые совершенно необходимы при использовании языка ассемблера.

Не все инструкции одинаково дорогие

В мирные, давно минувшие дни программирования на языке С Кернигана и Ритчи каждая инструкция была примерно такой же дорогой, как и любая другая. Вызов функции может содержать вычисления произвольной сложности. Однако оператор присваивания в общем случае копирует нечто, помещающееся в машинном регистре, во что-то другое, что может хранить содержимое регистра компьютера. Таким образом, инструкция

```
int i, j;  
...  
i = j;
```

копирует 2 или 4 байта из j в i . Объявление может быть `int`, `float` или `struct big_struct*`, но инструкция присваивания все равно выполняет одно и то же количество работы.

В настоящее время это не так. В C++ присвоение одного `int` другому представляет собой точно такой же объем работы, как и в случае соответствующей инструкции С. Но инструкция наподобие `BigInstance i = OtherObject;` может копировать целые структуры. Более того, этот вид присваивания вызывает конструктор `BigInstance`, который может скрывать произвольно сложный механизм. Конструктор вызывается также для каждого выражения, переданного функции в качестве формального аргумента, и вновь, когда функция возвращает значение. Арифметические операторы и операторы сравнения также могут быть перегружены, так что выражение `A = B * C;` может умножать n -мерные матрицы, а `if (x < y) ...` может сравнивать два пути через ориентированный граф произвольной сложности. Значение этого для оптимизации заключается в том, что *некоторые инструкции скрывают большие количества вычислений. Вид инструкции ничего не говорит об ее стоимости.*

Разработчики, которые начинали изучение языков программирования с C++, могут не увидеть в этом ничего удивительного; но для тех, кто начинал с изучения С, их инстинкты могут привести к катастрофическим заблуждениям.

Инструкции выполняются не по порядку

Программы на C++ ведут себя так, как если бы они выполнялись в порядке, указанном инструкциями управления потоком выполнения C++. Хитрая оговорка “как если бы” в предыдущем предложении является главной, на которой построено множество оптимизаций и трюков современного компьютерного оборудования.

За кулисами компилятор может — и зачастую так и делает — переупорядочивать инструкции для повышения производительности. Но компилятор знает, что переменная должна содержать последний результат вычисления, присвоенный до проверки или присваивания его другой переменной. Современные микропроцессоры также могут выбрать команды для выполнения не по порядку их следования, но они содержат логические схемы, которые гарантируют выполнение записи в память до последующего чтения того же места. Логика управления памятью микропроцессора может даже выбрать задержку записи в память для оптимального использования шины памяти. Однако контроллер памяти знает, что именно в настоящее время находится в

процессе записи от исполняющего устройства через кеш-память в основную память, и гарантирует, что если будет считываться тот же адрес, то будет получено корректное значение, даже если его запись выполнена еще не до конца.

Параллелизм усложняет эту картину. Программы на C++ компилируются без знания о других потоках, которые могут выполняться одновременно с данным. Компилятор C++ не знает, какие переменные, если таковые имеются, совместно используются потоками. Совокупный эффект изменения порядка инструкций компилятором и компьютером и задержка записи в основную память разрушают иллюзию выполнения инструкций в указываемом программой порядке, если программа содержит параллельные потоки, которые совместно используют данные. Разработчик должен добавить явную синхронизацию кода многопоточных программ для гарантии получения согласованного предсказуемого поведения. *Синхронизация кода снижает степень параллелизма, получаемого при совместном использовании данных параллельными потоками.*

Резюме

- *Доступ к памяти доминирует над другими действиями процессора.*
- *Невыровненный доступ требует в два раза больше времени, чем когда все байты находятся в одном и том же слове.*
- *Доступ к интенсивно используемым ячейкам памяти можно получить быстрее, чем к используемым реже.*
- *Обращение к соседним ячейкам памяти выполняется быстрее, чем к памяти в отдаленных местах.*
- *Из-за кеширования определенная функция, работающая в контексте всей программы, может выполняться медленнее, чем та же функция в тестовой программе.*
- *Обращение к совместно используемым потоками выполнения данным выполняется гораздо медленнее, чем к не разделяемым данным.*
- *Вычисления осуществляются быстрее принятий решений.*
- *Каждая программа конкурирует с другими программами за ресурсы компьютера.*
- *Если программа должна запускаться при запуске системы или в периоды пиковой нагрузки, ее производительность должна измеряться под нагрузкой.*
- *Каждое присваивание, инициализация аргумента функции и возврат из функции вызывает конструктор — функцию, которая может скрывать произвольное большое количество кода.*
- *Некоторые инструкции скрывают большие количества вычислений. Вид инструкции ничего не говорит об ее стоимости.*
- *Синхронизация кода снижает степень параллелизма, получаемого при совместном использовании данных параллельными потоками.*