

Ответы на контрольные вопросы

1. В чем основные отличия алгоритмов сортировки пузырьком, выбором и вставкой?

Все три алгоритма простые и имеют сложность $O(n^2)$, но работают по-разному.

Сортировка пузырьком: Сравнивает соседние элементы и меняет их местами если они не в порядке. Самые большие элементы постепенно "всплывают" к концу массива как пузырьки. Делает много обменов элементов - это самый медленный из трех алгоритмов. Проходит по массиву несколько раз пока он не станет полностью отсортированным.

Пример: [5, 2] -> видим что $5 > 2$, меняем -> [2, 5]. Потом сравниваем следующую пару.

Сортировка выбором: На каждом шаге находит минимальный элемент в оставшейся части массива и ставит его на текущую позицию. Делает мало обменов (ровно n обменов), но много сравнений. Работает чуть быстрее пузырька. Всегда делает одинаковое количество операций независимо от входных данных.

Пример: в массиве [5, 2, 8, 1] находим минимум (1), ставим в начало -> [1, 2, 8, 5]. Теперь ищем минимум среди [2, 8, 5].

Сортировка вставкой: Строит отсортированную часть массива постепенно, вставляя каждый новый элемент на правильное место. Работает как сортировка карт в руке. На почти отсортированных данных работает очень быстро - почти за $O(n)$. Из трех алгоритмов самая эффективная на практике.

Пример: [5 | 2, 8, 1] - пять уже на месте, берем двойку, вставляем перед пятеркой -> [2, 5 | 8, 1].

Главное отличие в подходе:

- Пузырек: меняй соседей пока не будет порядка
- Выбор: найди минимум и поставь на место
- Вставка: бери элементы и вставляй в отсортированную часть

На практике: Вставка обычно быстрее всех на реальных данных. Выбор делает меньше обменов чем пузырек. Пузырек самый медленный, его используют только в учебных целях.

2. Почему параллельная реализация сортировки вставкой сложнее для выполнения с использованием OpenMP?

Сортировка вставкой по своей природе последовательный алгоритм. Каждый шаг зависит от результата всех предыдущих шагов.

Проблема зависимости данных: Когда мы вставляем третий элемент, нам нужно чтобы первые два уже были отсортированы. Когда вставляем четвертый - первые три должны быть на местах. То есть каждый шаг ждет завершения предыдущего. Нельзя вставить пятый элемент пока не вставлен четвертый.

Представь очередь в магазине. Третий человек не может пройти к кассе пока не обслужили второго. А второй ждет первого. Нельзя обслужить всех параллельно.

Почему это проблема для OpenMP: OpenMP хорошо работает когда можно разделить работу на независимые куски. Например при поиске минимума каждый поток может искать в своей части массива, потом просто выбираем минимум из результатов. Тут все просто.

Но в сортировке вставкой нельзя вставлять несколько элементов одновременно - они будут мешать друг другу и результат будет неправильным.

Что можно сделать: Мы использовали гибридный подход:

1. Разбиваем массив на куски
2. Каждый поток сортирует свой кусок независимо (это параллельно)
3. Потом сливаем отсортированные куски (это последовательно)

Но это уже не чистая сортировка вставкой, а скорее комбинация вставки и слияния.

Результат: Из-за этих проблем параллельная сортировка вставкой дает меньше ускорения чем можно было бы ожидать. Большая часть работы все равно выполняется последовательно.

3. Какие директивы OpenMP были использованы для параллельной реализации алгоритмов?

В коде использовались несколько директив OpenMP, каждая для своей цели.

#pragma omp parallel for Самая часто используемая директива. Делит итерации цикла между потоками. Используется в пузырьке для параллельного сравнения пар элементов.

```
#pragma omp parallel for
for (int i = 0; i < n - 1; i += 2) {
    // каждый поток обрабатывает свои пары
}
```

Автоматически создает потоки, делит работу, потом ждет завершения всех потоков.

#pragma omp parallel Создает команду потоков. Все что внутри блока выполняется всеми потоками. Используется в сортировке выбором когда нужен более тонкий контроль.

```
#pragma omp parallel
{
    // каждый поток выполняет этот код
}
```

#pragma omp for Используется внутри parallel блока для распределения итераций цикла. Отличается от parallel for тем что не создает новые потоки, а использует уже существующие.

```
#pragma omp parallel
{
    #pragma omp for
    for (int j = i + 1; j < n; j++) {
        // работа делится между потоками
    }
}
```

#pragma omp critical Критическая секция - участок кода где может находиться только один поток. Используется в сортировке выбором для сравнения локальных минимумов.

```
#pragma omp critical
{
    if (localMinValue < minValue) {
        minValue = localMinValue;
    }
}
```

Если один поток внутри критической секции, остальные ждут снаружи.

Почему эти директивы:

- parallel for - когда нужно просто распараллелить цикл
- parallel + for - когда нужны локальные переменные для каждого потока
- critical - когда нужно защитить общую переменную от одновременного доступа

4. Какие преимущества и недостатки параллельной реализации алгоритмов сортировки на CPU?

Преимущества:

1. Ускорение на больших данных Если массив большой и процессор имеет несколько ядер, можем получить ускорение в 1.5-3 раза. На 4 ядрах теоретически можем ускорить в 4 раза, но на практике меньше из-за накладных расходов.

2. Лучшее использование ресурсов Современные процессоры имеют 4-16 ядер. Последовательная программа использует только одно ядро, остальные пропстаивают. Параллельная версия загружает все ядра работой.

3. Масштабируемость На процессоре с большим количеством ядер получим большее ускорение. Код автоматически адаптируется под количество доступных ядер.

4. Простота реализации с OpenMP Не нужно вручную создавать потоки и управлять ими. Достаточно добавить несколько директив в существующий код.

Недостатки:

1. Накладные расходы Создание потоков, их синхронизация, объединение результатов требуют времени. На маленьких массивах (меньше 1000 элементов) эти расходы могут быть больше чем выигрыш от параллелизма.

2. Сложность отладки Баги в параллельном коде сложнее найти. Гонки данных (race conditions) могут проявляться не всегда, результат может быть недетерминированным.

3. Ограничено ускорение для $O(n^2)$ алгоритмов Простые алгоритмы сортировки плохо параллелятся из-за зависимостей между шагами. Даже с 8 ядрами редко получаем ускорение больше 3x.

4. Нужна синхронизация Критические секции и барьеры синхронизации замедляют программу. Потоки проводят время в ожидании друг друга вместо полезной работы.

5. Неэффективность на частично отсортированных данных Последовательная сортировка вставкой на почти отсортированных данных работает очень быстро. Параллельная версия теряет это преимущество из-за накладных расходов.

6. Проблемы с кешем процессора Когда несколько потоков работают с одними данными, возникает false sharing - потоки постоянно инвалидируют кеш друг друга. Это сильно замедляет работу.

Вывод: Параллелизация простых сортировок дает небольшой выигрыш и имеет много ограничений. Для реального ускорения сортировки лучше использовать более сложные алгоритмы (quick sort, merge sort) которые лучше параллелятся.

5. Как можно измерить производительность программы в C++?

В C++ есть несколько способов измерить время выполнения кода. Самый современный и точный - библиотека chrono.

Использование библиотеки chrono:

```
#include <chrono>

// запоминаем время начала
auto start = chrono::high_resolution_clock::now();

// код который измеряем
sortArray(arr);

// запоминаем время конца
auto end = chrono::high_resolution_clock::now();

// вычисляем разницу
auto duration = chrono::duration_cast<chrono::milliseconds>(end - start);

// выводим результат
cout << "Время: " << duration.count() << " мс" << endl;
```

Почему chrono:

- Высокая точность (до наносекунд)
- Кроссплатформенность (работает везде одинаково)
- Не зависит от системных часов
- Современный C++ стандарт

Единицы измерения:

- **nanoseconds** - наносекунды (10^{-9} секунды) - для очень быстрых операций
- **microseconds** - микросекунды (10^{-6} секунды) - для быстрых функций
- **milliseconds** - миллисекунды (10^{-3} секунды) - для обычных операций
- **seconds** - секунды - для долгих вычислений

Важные моменты при измерении:

1. Прогрев (warmup) Первый запуск может быть медленнее из-за холодного кеша.
Лучше запустить функцию несколько раз перед измерением.

2. Несколько измерений Время может колебаться из-за других процессов в системе.
Лучше измерить 5-10 раз и взять среднее.

```
double totalTime = 0;
for (int i = 0; i < 10; i++) {
    auto start = chrono::high_resolution_clock::now();
    sortArray(arr);
    auto end = chrono::high_resolution_clock::now();
    auto duration = chrono::duration_cast<chrono::milliseconds>(end - start);
    totalTime += duration.count();
}
```

```
auto end = chrono::high_resolution_clock::now();
auto duration = chrono::duration_cast<chrono::milliseconds>(end - start);
totalTime += duration.count();
}
double avgTime = totalTime / 10;
```

3. Компиляция с оптимизацией Без оптимизации код работает медленнее.

Компилировать надо с флагом -O2 или -O3:

```
g++ -O3 -fopenmp program.cpp -o program
```

4. Отключить другие программы Браузер, видео, игры влияют на время. Для точных измерений закрыть все лишнее.

Старые способы (не рекомендуются):

clock() из <ctime>:

```
clock_t start = clock();
sortArray(arr);
clock_t end = clock();
double time = (double)(end - start) / CLOCKS_PER_SEC;
```

Проблема: низкая точность, измеряет процессорное время а не реальное.

gettimeofday(): Только в Unix системах, устаревший подход.

Вывод: Используйте chrono::high_resolution_clock для точных измерений. Делайте несколько измерений и берите среднее. Компилируйте с оптимизацией.

6. Как изменяется производительность сортировок при увеличении числа потоков?

Зависимость производительности от числа потоков нелинейная и имеет предел.

Теоретическая картина: Если у нас 1 поток и мы добавляем второй, работа должна выполняться в 2 раза быстрее. С 4 потоками - в 4 раза быстрее. Но на практике все сложнее.

Реальная картина для простых сортировок:

1-2 потока: Переход от 1 к 2 потокам дает заметное ускорение - примерно в 1.5 раза. Это лучшее соотношение, так как накладных расходов еще немного.

2-4 потока: Ускорение продолжается но медленнее. С 4 потоками общее ускорение относительно 1 потока обычно 2-2.5 раза вместо теоретических 4 раз.

4-8 потоков: Прирост становится минимальным. Ускорение может быть 2.5-3 раза максимум. Дальнейшее увеличение потоков почти не помогает.

Больше 8 потоков: Производительность может даже падать из-за избыточных накладных расходов на управление потоками.

Почему не получаем линейное ускорение:

1. Закон Амдала Часть алгоритма нельзя распараллелить (последовательная часть). Эта часть ограничивает максимальное ускорение. Для сортировок $O(n^2)$ последовательная часть значительная.

Формула: Ускорение = $1 / (S + (1-S)/N)$ где S - доля последовательного кода, N - число потоков.

Если 50% кода последовательный, максимальное ускорение = 2x независимо от числа потоков.

2. Синхронизация Чем больше потоков, тем больше времени тратится на синхронизацию. Критические секции становятся узким местом - потоки стоят в очереди.

3. False sharing Когда потоки работают с данными в одной строке кеша, они постоянно инвалидируют кеш друг друга. С ростом числа потоков проблема усиливается.

4. Перегрузка планировщика Операционная система тратит время на переключение между потоками. Если потоков больше чем ядер процессора, много времени уходит на переключение контекста.

Практические наблюдения:

Для сортировки пузырьком:

- 1 поток: базовое время
- 2 потока: ускорение 1.3-1.5x
- 4 потока: ускорение 1.8-2.2x
- 8 потоков: ускорение 2-2.5x

Для сортировки выбором:

- Похожие цифры, может быть чуть лучше

Для сортировки вставкой (гибрид):

- Немного лучше так как этап сортировки кусков полностью параллелен

Оптимальное число потоков: Обычно лучший результат при количестве потоков равном числу физических ядер процессора. На процессоре с 4 ядрами оптимально 4 потока.

Вывод: Увеличение потоков дает убывающую отдачу. После 4 потоков прирост минимален для простых сортировок. Это из-за последовательных участков кода и накладных расходов на синхронизацию.

7. В каких ситуациях параллельная сортировка может быть менее эффективной чем последовательная?

Есть несколько случаев когда параллелизм вредит производительности.

1. Маленькие массивы (меньше 1000 элементов)

Накладные расходы на создание потоков и синхронизацию занимают больше времени чем сама сортировка.

Пример: сортировка 100 элементов занимает 0.01 мс. Создание 4 потоков и их синхронизация - 0.1 мс. Итого параллельная версия в 10 раз медленнее.

Правило: если массив меньше 1000-5000 элементов для $O(n^2)$ сортировок - используй последовательную версию.

2. Почти отсортированные данные

Последовательная сортировка вставкой на почти отсортированных данных работает почти за $O(n)$ - очень быстро. Она просто проверяет что элементы уже на месте.

Параллельная версия не может использовать это преимущество. Она все равно делит массив на куски, сортирует их, потом сливает. Все эти операции требуют времени даже если массив уже отсортирован.

Результат: на почти отсортированных данных последовательная версия может быть в 5-10 раз быстрее параллельной.

3. Мало доступных ядер процессора

Если у процессора только 2 ядра или другие программы используют ядра, параллелизм дает мало выигрыша. Создание 4 потоков на 2 ядрах приводит к постоянному переключению контекста.

4. Много других активных процессов

Если система загружена (открыт браузер, идет компиляция, работает антивирус), ядра процессора заняты. Потоки твоей программы получают меньше процессорного времени и постоянно вытесняются.

Результат: последовательная версия работает стабильнее в таких условиях.

5. Ограниченнaя память

Некоторые параллельные версии (например наша сортировка вставкой) требуют дополнительную память для временных массивов. Если памяти мало, это вызывает свопинг (выгрузку на диск) что катастрофически замедляет работу.

6. Данные с плохой локальностью

Если элементы массива разбросаны по памяти (например массив указателей), параллельный доступ вызывает много cache miss (промахов кеша). Процессор тратит время на чтение из медленной оперативной памяти.

Последовательная версия лучше использует кеш благодаря предсказанию обращений к памяти.

7. Отладочная версия без оптимизации

Если компилировать без флага `-O2/-O3`, параллельная версия будет намного медленнее. Компилятор не оптимизирует код, добавляет много проверок, накладные расходы OpenMP становятся очень большими.

8. Специфика алгоритма

Для сортировки вставкой параллелизм почти бесполезен из-за зависимостей данных. Иногда последовательная версия с хорошо написанным кодом и оптимизациями компилятора работает быстрее кривой параллельной версии.

Вывод:

Параллелизм - не серебряная пуля. Он эффективен только когда:

- Данных много (больше 10000 элементов)
- Данные случайные или отсортированы плохо
- Процессор имеет несколько свободных ядер
- Система не перегружена
- Код скомпилирован с оптимизациями

В остальных случаях лучше использовать последовательную версию или вообще другой алгоритм (quick sort, merge sort с $O(n \log n)$).