

ГЛАВА 6

ПАРАЛЛЕЛЬНЫЕ МЕТОДЫ УМНОЖЕНИЯ МАТРИЦЫ НА ВЕКТОР

6.1. Введение

Матрицы и матричные операции широко используются при математическом моделировании самых разнообразных процессов, явлений и систем. Матричные вычисления составляют основу многих научных и инженерных расчетов – среди областей приложений могут быть указаны вычислительная математика, физика, экономика и др.

С учетом значимости эффективного выполнения матричных расчетов многие стандартные библиотеки программ содержат процедуры для различных матричных операций. Объем программного обеспечения для обработки матриц постоянно увеличивается – разрабатываются новые экономные структуры хранения для матриц специального типа (треугольных, ленточных, разреженных и т. п.), создаются различные высокоэффективные машинно-зависимые реализации алгоритмов, проводятся теоретические исследования для поиска более быстрых методов матричных вычислений.

Являясь вычислительно-трудоемкими, матричные вычисления представляют собой классическую область применения параллельных вычислений. С одной стороны, использование высокопроизводительных многопроцессорных систем позволяет существенно повысить сложность решаемых задач. С другой стороны, в силу своей достаточно простой формулировки матричные операции дают прекрасную возможность для демонстрации многих приемов и методов параллельного программирования.

В данной главе обсуждаются методы параллельных вычислений для операции матрично-векторного умножения, в следующей главе будет рассмотрена операция перемножения матриц. Важный вид матричных вычислений – решение систем линейных уравнений – представлен в главе 8. Общий для всех перечисленных задач вопрос разделения обрабатываемых матриц между параллельно работающими потоками рассматривается в разделе 6.2.

При изложении следующего материала будем полагать, что рассматриваемые матрицы являются *плотными* (*dense*), в которых число нулевых элементов является незначительным по сравнению с общим количеством элементов матриц.

6.2. Принципы распараллеливания

Для многих методов матричных вычислений характерным является повторение одних и тех же вычислительных действий для разных элементов матриц. Данный момент свидетельствует о наличии *параллелизма по данным* при выполнении матричных расчетов и, как результат, распараллеливание матричных операций сводится в большинстве случаев к разделению обрабатываемых матриц между потоками. Выбор способа разделения матриц приводит к определению конкретного метода параллельных вычислений; существование разных схем распределения данных порождает целый ряд *параллельных алгоритмов матричных вычислений*.

Наиболее общие и широко используемые способы разделения матриц состоят в разбиении данных на *полосы* (по вертикали или горизонтали) или на прямоугольные фрагменты (*блоки*).

1. Ленточное разбиение матрицы. При *ленточном* (*block-striped*) разбиении каждому потоку выделяется то или иное подмножество строк (*rowwise* или *горизонтальное разбиение*) или столбцов (*columnwise* или *вертикальное разбиение*) матрицы (рис. 6.1). Разделение строк и столбцов на полосы в большинстве случаев происходит на *непрерывной* (*последовательной*) основе. При таком подходе для горизонтального разбиения по строкам, например, матрица A представляется в виде (см. рис. 6.1)

$$A = (A_0, A_1, \dots, A_{p-1})^T, A_i = (a_{i_0}, a_{i_1}, \dots, a_{i_{k-1}}), i_j = ik + j, 0 \leq j < k, k = m/p, \quad (6.1)$$

где $a_i = (a_{i1}, a_{i2}, \dots, a_{in})$, $0 \leq i < m$, есть i -я строка матрицы A (предполагается, что количество строк m кратно числу вычислительных элементов (процессоров и/или ядер) p , т. е. $m = k \cdot p$). Во всех алгоритмах матричного умножения и умножения матрицы на вектор, которые будут рассмотрены нами в этом и следующем разделах, используется разделение данных на непрерывной основе.

Другой возможный подход к формированию полос состоит в применении той или иной схемы *чередования* (*циклическости*) строк или столбцов. Как правило, для чередования используется число потоков p – в этом случае при горизонтальном разбиении матрица A принимает вид

$$A = (A_0, A_1, \dots, A_{p-1})^T, A_i = (a_{i_0}, a_{i_1}, \dots, a_{i_{k-1}}), i_j = i + jp, 0 \leq j < k, k = m/p. \quad (6.2)$$

Циклическая схема формирования полос может оказаться полезной для лучшей балансировки вычислительной нагрузки вычислительных элементов (например, при решении системы линейных уравнений с использованием метода Гаусса – см. главу 8).

2. Блочное разбиение матрицы. При *блочном* (*chessboard block*) разделении матрица делится на прямоугольные наборы элементов – при этом,

как правило, используется разделение на непрерывной основе. Пусть количество потоков составляет $p = s \cdot q$, количество строк матрицы является кратным s , а количество столбцов – кратным q , то есть $m = k \cdot s$ и $n = l \cdot q$. Представим исходную матрицу A в виде набора прямоугольных блоков следующим образом:

$$A = \begin{pmatrix} A_{00} & A_{02} & \dots A_{0q-1} \\ & \dots & \\ A_{s-11} & A_{s-12} & \dots A_{s-1q-1} \end{pmatrix},$$

где A_{ij} – блок матрицы, состоящий из элементов:

$$A_{ij} = \begin{pmatrix} a_{i_0j_0} & a_{i_0j_1} & \dots a_{i_0j_{l-1}} \\ & \dots & \\ a_{i_{k-1}j_0} & a_{i_{k-1}j_1} & a_{i_{k-1}j_{l-1}} \end{pmatrix},$$

$$i_v = ik + v, 0 \leq v < k, k = m/s, j_u = jl + u, 0 \leq u \leq l, l = n/q. \quad (6.3)$$

При таком подходе часто оказывается полезным, чтобы топология вычислительной системы имела (по крайней мере, на логическом уровне) вид решетки из s строк и q столбцов. В этом случае при разделении данных на непрерывной основе вычисления в большинстве случаев могут быть организованы таким образом, чтобы вычислительные элементы, соседние в структуре решетки, обрабатывали смежные блоки исходной матрицы. Следует отметить, что и для блочной схемы может быть применено циклическое чередование строк и столбцов.

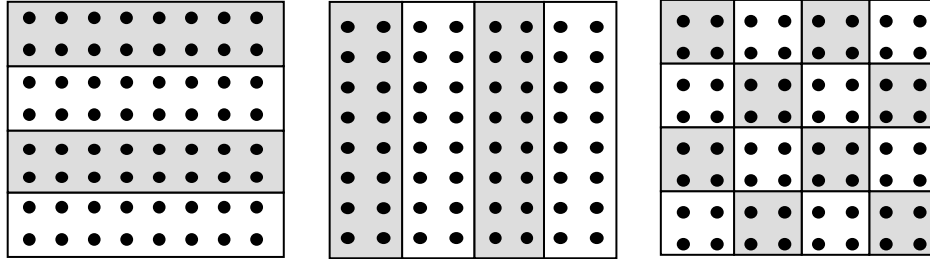


Рис. 6.1. Способы распределения элементов матрицы между потоками

В данной главе рассматриваются три параллельных алгоритма для умножения квадратной матрицы на вектор. Каждый подход основан на разном типе распределения исходных данных (элементов матрицы и вектора) между потоками. Для каждого рассматриваемого алгоритма проводится

теоретическая и экспериментальная оценка эффективности получаемых параллельных вычислений для определения наилучшего способа разделения данных.

6.3. Постановка задачи

В результате умножения матрицы A размера $m \times n$ и вектора b , состоящего из n элементов, получается вектор c размера m , каждый i -й элемент которого есть результат скалярного умножения i -й строки матрицы A (обозначим эту строчку a_i) и вектора b .

$$c_i = a_i, b = \sum_{j=1}^n a_{ij} b_j, 0 \leq i < m. \quad (6.4)$$

Тем самым, получение результирующего вектора c предполагает повторение m однотипных операций по умножению строк матрицы A и вектора b . Каждая такая операция включает умножение элементов строки матрицы и вектора b (n операций) и последующее суммирование полученных произведений ($n-1$ операций). Общее количество необходимых скалярных операций есть величина

$$T_1 = m \cdot (n - 1)$$

6.4. Последовательный алгоритм

Последовательный алгоритм умножения матрицы на вектор может быть представлен следующим образом.

```
// Программа 6.1
// Умножение матрицы на вектор
// (последовательный алгоритм)
for (i = 0; i < m; i++){
    c[i] = 0;
    for (j = 0; j < n; j++){
        c[i] += A[i][j]*b[j]
    }
}
```

Матрично-векторное умножение – это последовательность вычисления скалярных произведений. Поскольку каждое вычисление скалярного произведения векторов длины n требует выполнения n операций умножения и $n-1$ операций сложения, его трудоемкость порядка $O(n)$. Для выполнения матрично-векторного умножения необходимо выполнить m операций вы-

числения скалярного произведения; таким образом, алгоритм имеет трудоемкость порядка $O(mn)$.

При дальнейшем изложении материала для снижения сложности и упрощения получаемых соотношений будем предполагать, что матрица A является квадратной, т. е. $m = n$.

Представим возможный вариант последовательной программы умножения матрицы на вектор.

1. Главная функция программы. Программа реализует логику работы алгоритма, последовательно вызывает необходимые подпрограммы.

```
// Программа 6.2
// Умножение матрицы на вектор
// (последовательный алгоритм)
void main(int argc, char* argv[]) {
    double* pMatrix; // Исходная матрица
    double* pVector; // Исходный вектор
    double* pResult; // Вектор-результат
    int Size;        // Размер матрицы и векторов

    // Инициализация данных
    ProcessInitialization(pMatrix, pVector, pResult, Size);

    // Матрично-векторное умножение
    SerialResultCalculation(pMatrix, pVector, pResult, Size);

    // Завершение вычислений
    ProcessTermination(pMatrix, pVector, pResult);
}
```

2. Функция ProcessInitialization. Эта функция определяет размер и элементы для матрицы A и вектора b . Значения для матрицы A и вектора b определяются в функции *RandomDataInitialization*.

```
// Функция выделения памяти и инициализации данных
void ProcessInitialization (double* &pMatrix,
    double* &pVector, double* &pResult, int &Size) {
    int i;

    do {
        printf("\nВведите размер матрицы и векторов: ");
        scanf("%d", &Size);
        if (Size <= 0) {
            printf("Размер должен быть больше 0! \n ");
        }
    } while (Size <= 0);
}
```

```

    }
} while (Size <= 0);

pMatrix = new double [Size*Size];
pVector = new double [Size];
pResult = new double [Size];
for (i=0; i<Size; i++)
    pResult[i] = 0;
RandomDataInitialization(pMatrix, pVector, Size);
}

```

Реализация функции *RandomDataInitialization* предлагается для самостоятельного выполнения. Исходные данные могут быть введены с клавиатуры, прочитаны из файла или сгенерированы при помощи датчика случайных чисел.

3. Функция *SerialResultCalculation*. Данная функция производит умножение матрицы на вектор.

```

// Функция матрично-векторного умножения
void SerialResultCalculation(double* pMatrix,
    double* pVector, double* pResult, int Size) {
    int i, j;
    for (i=0; i<Size; i++) {
        for (j=0; j<Size; j++)
            pResult[i] += pMatrix[i*Size+j]*pVector[j];
    }
}

```

Разработку функции *ProcessTermination* также предлагается выполнить самостоятельно.

Рассмотрим возможные способы организации параллельных вычислений в задаче умножении матрицы на вектор на многопроцессорных вычислительных системах с общей памятью. Как и ранее, ради общности излагаемого учебного материала при упоминании одновременно и мультипроцессоров, и многоядерных процессоров, для обозначения одного вычислительного устройства будет использоваться понятие *вычислительного элемента (ВЭ)*.

6.5. Умножение матрицы на вектор при разделении данных по строкам

Рассмотрим в качестве первого примера организации параллельных матричных вычислений алгоритм умножения матрицы на вектор, основанный на представлении матрицы непрерывными наборами (горизонтальными

ми полосами) строк. При таком способе разделения данных в качестве базовой подзадачи может быть выбрана операция скалярного умножения одной строки матрицы на вектор. После завершения вычислений каждая базовая подзадача определяет один из элементов вектора результата s .

6.5.1. Выделение информационных зависимостей

В общем виде схема информационного взаимодействия подзадач в ходе выполняемых вычислений показана на рис. 6.2.

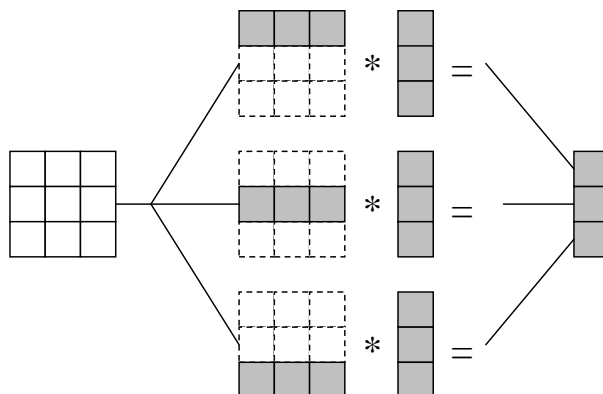


Рис. 6.2. Организация вычислений при выполнении параллельного алгоритма умножения матрицы на вектор, основанного на разделении матрицы по строкам

6.5.2. Масштабирование и распределение подзадач по вычислительным элементам

В процессе умножения плотной матрицы на вектор количество вычислительных операций для получения скалярного произведения одинаково для всех базовых подзадач. Поэтому когда число p вычислительных элементов меньше числа m базовых подзадач, следует объединить базовые подзадачи таким образом, чтобы каждый ВЭ выполнял несколько таких задач, соответствующих непрерывной последовательности строк матрицы A . В этом случае по окончании вычислений каждая базовая подзадача определяет набор элементов (блок) результирующего вектора s .

6.5.3. Программная реализация

Для того, чтобы разработать параллельную программу, реализующую описанный подход, при помощи технологии OpenMP, необходимо внести минимальные изменения в функцию умножения матрицы на вектор. Дос-

точно добавить одну директиву *parallel for* в функцию *SerialResultCalculation* (назовем новый вариант функции *ParallelResultCalculation*):

```
// Программа 6.3
// Функция параллельного матрично-векторного умножения
void ParallelResultCalculation(double* pMatrix,
    double* pVector, double* pResult, int Size) {
    int i, j;
    #pragma omp parallel for private (j)
    for (i=0; i<Size; i++) {
        for (j=0; j<Size; j++)
            pResult[i] += pMatrix[i*Size+j]*pVector[j];
    }
}
```

Данная функция производит умножение строк матрицы на вектор с использованием нескольких параллельных потоков. Каждый поток выполняет вычисления над несколькими соседними строками матрицы i , таким образом, получает блок результирующего вектора c .

Параллельные области в данной функции задаются директивой *parallel for*. Компилятор, поддерживающий технологию OpenMP, разделяет выполнение итераций цикла между несколькими потоками параллельной программы, количество которых обычно совпадает с числом вычислительных элементов (процессоров и/или ядер) в вычислительной системе. Параметр директивы *private* указывает необходимость создания отдельных копий для задаваемых переменных для каждого потока — эти копии могут использоваться в потоках независимо друг от друга.

В данной функции между потоками параллельной программы разделяются итерации внешнего цикла алгоритма умножения матрицы на вектор. В результате каждый поток выполняет непрерывную последовательность итераций цикла по переменной i и, таким образом, умножает горизонтальную полосу матрицы A на вектор b и вычисляет блок элементов результирующего вектора c .

6.5.4. Анализ эффективности

Алгоритм умножения матрицы на вектор, основанный на ленточном горизонтальном разбиении матрицы, обладает хорошей «локализацией вычислений», т. е. каждый поток параллельной программы использует только «свои» данные, и ему не требуются данные, которые в данный момент обрабатывает другой поток, нет обмена данными между потоками, не возникает необходимости синхронизации. Это означает, что в данной задаче практически нет накладных расходов на организацию па-

раллелизма (за исключением расходов на организацию и закрытие параллельной секции), и можно ожидать линейного ускорения. Однако, как будет далее видно из представленных результатов (см. табл. 6.3), ускорение, которое демонстрирует параллельный алгоритм умножения матрицы на вектор, основанный на ленточном горизонтальном разделении данных, далек от линейного. В чем же причина такого поведения параллельной программы?

Задача умножения матрицы на вектор обладает сравнительно невысокой вычислительной сложностью – трудоемкость алгоритма имеет порядок $O(n^2)$. Такой же порядок – $O(n^2)$ – имеет и объем данных, обрабатываемый алгоритмом умножения. Время решения задачи одним потоком складывается из времени, когда процессор непосредственно выполняет вычисления, и времени, которое тратится на чтение необходимых для вычислений данных из оперативной памяти в кэш память. При этом время, необходимое на чтение данных, может быть сопоставимо или даже превосходить время счета.

Проведем простой эксперимент: измерим время выполнения последовательной программы, которая суммирует все элементы матрицы, и сравним это время со временем выполнения умножения матрицы того же размера на вектор. Результаты вычислительных экспериментов приведены в табл. 6.1 и представлены на рис. 6.3.

При суммировании элементов матрицы на каждой итерации цикла выполняется простая операция сложения двух чисел. При умножении матрицы на вектор на каждой итерации цикла выполняются две операции: более сложная операция умножения и операция сложения. Но, несмотря на большую вычислительную сложность, время работы алгоритма умножения матрицы на вектор превосходит время выполнения сложения элементов матрицы всего на 16%. Этот эксперимент можно рассматривать, как подтверждение предположения о том, что значительная часть времени тратится на выборку необходимых данных из оперативной памяти в кэш процессора.

Таблица 6.1.
Сравнение времени выполнения алгоритма умножения матрицы на вектор и алгоритма суммирования всех элементов матрицы

Размер матрицы	Время выполнения умножения матрицы на вектор	Время выполнения суммирования элементов матрицы
1000	0,0076	0,0064
2000	0,0303	0,0260
3000	0,0687	0,0577

4000	0,1221	0,1023
5000	0,1909	0,1593
6000	0,2747	0,2288
7000	0,3741	0,3123
8000	0,4893	0,4082
9000	0,6186	0,5152
10000	0,7637	0,6364

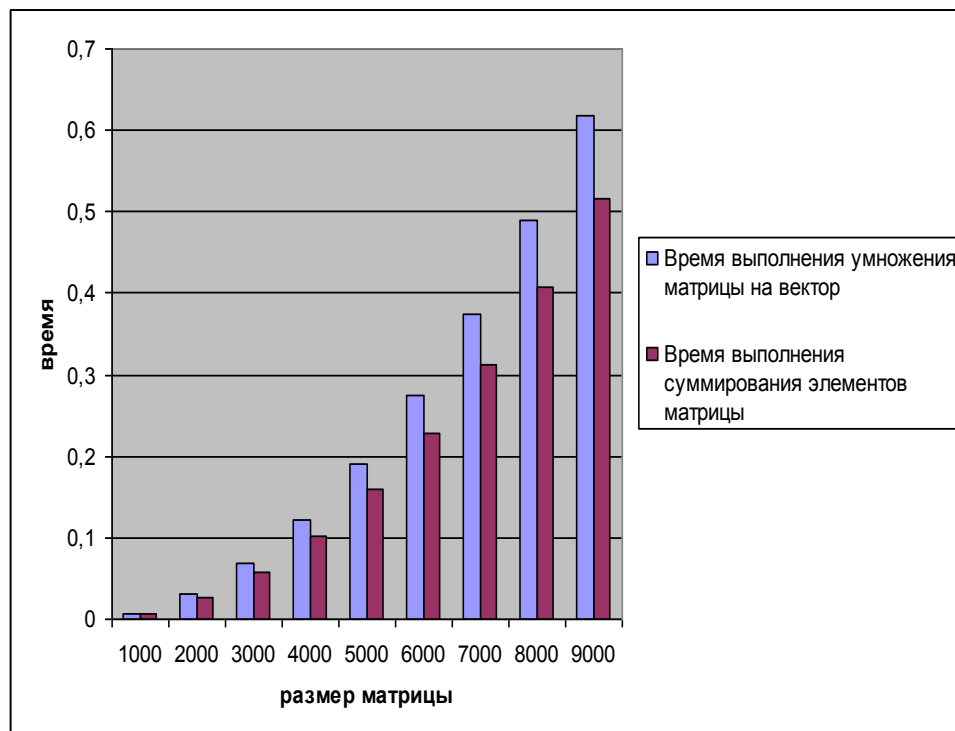


Рис. 6.3. Время выполнения умножения матрицы на вектор и суммирования всех элементов матрицы

Процесс выполнения последовательного алгоритма умножения матрицы на вектор можно представить в виде диаграммы (рис. 6.4).

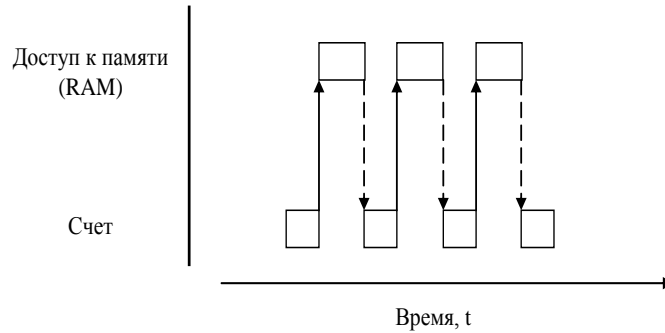


Рис. 6.4. Диаграмма состояний процесса выполнения последовательного алгоритма умножения матрицы на вектор

Следует отметить, что схема, представленная на рис. 6.4, является определенным упрощением процесса вычислений, реально выполняемых в компьютере, поскольку в современных процессорах на аппаратном уровне часто реализуются те или иные алгоритмы предсказания потребности данных для обработки, так что часть данных может перемещаться в кэш заранее (обеспечивая тем самым совмещение обработки и подкачки данных). Однако для предварительного анализа (тем более, для быстро выполняемых операций обработки) использование такой упрощенной схемы рис. 6.4. является вполне допустимой. Дополнительная информация по архитектуре современных процессоров может быть получена, например, в [35].

С учетом проведенного обсуждения, построим оценку времени выполнения матрично-векторного умножения. Итак, время выполнения последовательного алгоритма складывается из времени вычислений и времени доступа к памяти:

$$T_1 = T_{calc} + T_{mem} \quad (6.5)$$

Время выполнения вычислений можно определить как произведение количества выполненных операций N на время τ выполнения одной операции:

$$T_{calc} = n \cdot (2n - 1) \cdot \tau,$$

где n – порядок матрицы (для получения более простых зависимостей предполагается, что все выполняемые вычислительные операции являются одинаковыми и имеют одну и ту же длительность τ).

Время доступа к памяти можно получить как результат деления объема извлекаемых из памяти данных M (в данном случае $M = n^2$) на *пропуск-*

ную способность памяти¹⁾ β . Для точного построения модели следует учитывать, что данные загружаются в кэш память не побайтно, а целыми блоками – *кэш строками*. В вычислительной системе, на которой проводились все приводимые далее эксперименты, размер кэш-строки был равен 64 байтам. Кроме того, при построении оценок необходимо принимать во внимание, что выполнению операций доступа к памяти может предшествовать некоторая задержка (*латентность доступа к памяти*) α .

В предельном случае, каждое обращение к данным приводит к чтению из памяти отдельной кэш строки и, как результат, время доступа к памяти можно оценить как:

$$T_{mem} = n^2 \cdot (\alpha + 64/\beta) \quad (6.6)$$

Объединяя (6.5) и (6.6), получим суммарное время выполнения алгоритма:

$$T_1 = n \cdot 2n - 1 \cdot \tau + n^2 \cdot (\alpha + 64/\beta) \quad (6.7)$$

Следует отметить, что построенная оценка рассчитана на наихудший случай, когда любое обращение к данным будет приводить к необходимости доступа к памяти. На самом деле необходимые данные могут располагаться в кэш памяти и в этом случае необходимость доступа к оперативной памяти отпадает. Кэш память является существенно более быстрым видом памяти и использование данных из кэша приводит к заметному ускорению вычислений.

Потребность доступа к оперативной памяти возникает в моменты, когда необходимые данные отсутствуют в кэш памяти (такие ситуации носят наименование *кэш-промахов*). Тем самым, фактические затраты на работу с памятью во многом определяются соотношением использования кэш памяти и более медленной оперативной памяти. Если данные в основном берутся из кэш памяти, затраты на доступ к памяти становятся незначительными (кэш память работает практически на частоте процессора). Если частота кэш-промахов является большой, тогда затраты на доступ к памяти определяются быстродействием оперативной памяти.

Введем в разработанную ранее модель (6.7) величину γ , $0 \leq \gamma \leq 1$, для задания *частоты возникновения кэш-промахов*. Тогда оценка времени выполнения алгоритма матрично-векторного умножения принимает вид:

$$T_1 = n \cdot 2n - 1 \cdot \tau + \gamma \cdot n^2 \cdot (\alpha + 64/\beta) \quad (6.8)$$

Сведем воедино параметры построенной модели:

¹⁾ Пропускная способность памяти есть объем данных, который может быть извлечен из памяти за некоторую единицу времени

- n – размерность матрицы,
- α – латентность оперативной памяти,
- β – пропускная способность оперативной памяти,
- γ – частота кэш промахов,
- τ – время выполнения базовой вычислительной операции.

Рассмотрим способы оценки значений этих параметров.

Параметр n является исходной величиной и задается при постановке задачи.

Следующие два параметра α и β являются характеристиками используемой вычислительной системы и для их измерения существует ряд систем анализа производительности памяти. В приводимых далее экспериментах использовалась система RightMark Memory Analyzer [88], с помощью которой для задействованной в экспериментах вычислительной системе были получены оценки $\beta = 12,44$ Гб/с и $\alpha = 8,31$ нс.

Для оценки частоты кэш-промахов γ можно использовать систему Intel VTune [107]. Полезным средством для измерения данного параметра может оказаться и система Visual Performance System [24] (далее VPS), с помощью которой получались приводимые далее оценки.

Для того, чтобы оценить время τ одной операции, можно измерить время выполнения последовательного алгоритма умножения матрицы на вектор при малых объемах данных, таких, чтобы матрица и вектор полностью поместились в кэш памяти. Чтобы исключить необходимость выборки данных из оперативной памяти, перед началом вычислений можно заполнить матрицу и вектор случайными числами – выполнение этого действия гарантирует размещение данных в кэш памяти. Далее при решении задачи все время будет тратиться непосредственно на вычисления, т. к. нет необходимости загружать данные из оперативной памяти. Поделив полученное время на количество выполненных операций, можно получить время τ выполнения базовой вычислительной для алгоритма операции. Для вычислительной системы, которая далее использовалась для проведения экспериментов, было получено значение τ , равное 3,78 нс.

Получив оценки значений параметров, можно сравнить соответствие времен реального выполнения алгоритма матрично-векторного умножения и величин, получаемых с использованием разработанной модели (6.8) ²⁾.

Частота кэш-промахов, измеренная с помощью системы VPS, оказалась равной 0,0033.

²⁾ Описание использованной вычислительной системы приведено в п. 6.5.5.

Таблица 6.2

Сравнение экспериментального и теоретического времени выполнения последовательного алгоритма умножения матрицы на вектор

Размер матриц	Эксперимент	Модель
1000	0,0076	0,0076
2000	0,0303	0,0304
3000	0,0687	0,0684
4000	0,1221	0,1216
5000	0,1909	0,1901
6000	0,2747	0,2737
7000	0,3741	0,3725
8000	0,4893	0,4866
9000	0,6186	0,6158
10000	0,7637	0,7603

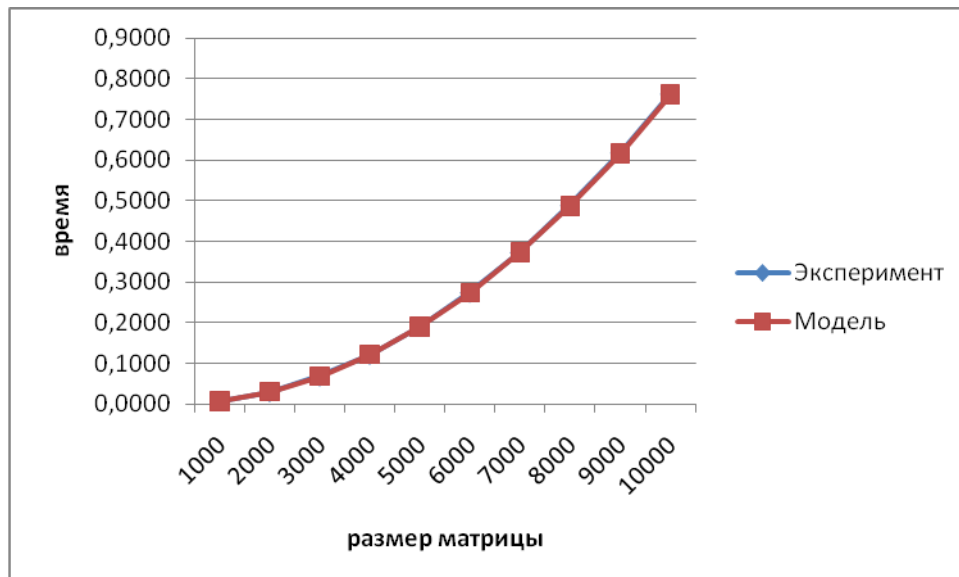


Рис. 6.5. График зависимости экспериментального и теоретического времен выполнения последовательного алгоритма от объема исходных данных (ленточное разбиение матрицы по строкам)

Как следует из приведенных результатов, разработанная модель с достаточно высокой степенью точности позволяет оценить время выполнения последовательного алгоритма матрично-векторного умножения.

Рассмотрим теперь вопрос оценки времени выполнения параллельного алгоритма матрично-векторного умножения. В многоядерных процессорах Intel архитектуры Core 2 Duo ядра процессоров разделяют общий канал доступа к памяти. Тем самым, несмотря на то, что вычисления могут выполняться ядрами параллельно, доступ к памяти осуществляется строго последовательно. На рис. 6.6 представлена возможная диаграмма состояний потоков, выполняющих параллельный алгоритм умножения матрицы на вектор, при условии, что эти потоки запущены на ядрах одного двухъядерного процессора:

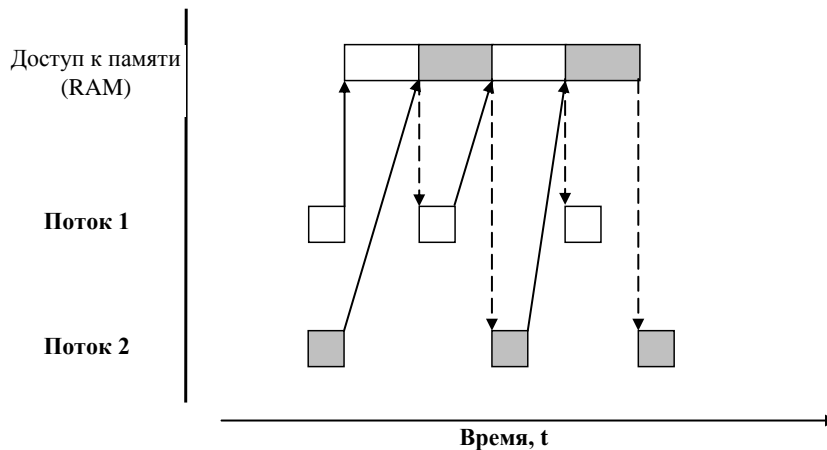


Рис. 6.6. Диаграмма состояний потоков при выполнении параллельного алгоритма умножения матрицы на вектор

Следовательно, время T_p выполнения параллельного алгоритма на компьютерной системе, имеющей p вычислительных элементов, с использованием p потоков и при описанной выше схеме доступа к памяти может быть оценено при помощи следующего соотношения:

$$T_p = \frac{n \cdot 2n - 1}{p} \cdot \tau + \gamma \cdot n^2 \cdot (\alpha + 64/\beta) \quad (6.9)$$

6.5.5. Результаты вычислительных экспериментов

Рассмотрим результаты вычислительных экспериментов, выполненных для оценки эффективности параллельного алгоритма умножения матрицы на вектор. Эксперименты проводились на двухпроцессорном вычислительном узле на базе четырехъядерных процессоров Intel Xeon E5320, 1.86 ГГц, 4 Гб RAM под управлением операционной системы Microsoft Windows HPC Server 2008. Разработка программ проводилась в среде Microsoft Visual Studio 2008, для компиляции использовался Intel C++ Compiler 10.0 for Windows.

Таблица 6.3.
Результаты вычислительных экспериментов для параллельного
алгоритма умножения матрицы на вектор при ленточной схеме
разделении данных по строкам

Размер матрицы	Последовательный алгоритм	Параллельный алгоритм			
		2 потока		4 потока	
		Время	Ускорение	Время	Ускорение
1000	0,0076	0,0038	2,0224	0,0027	2,8148
2000	0,0303	0,0150	2,0233	0,0100	3,0315
3000	0,0688	0,0341	2,0167	0,0227	3,0301
4000	0,1222	0,0606	2,0163	0,0412	2,9654
5000	0,1909	0,0947	2,0164	0,0626	3,0500
6000	0,2748	0,1356	2,0263	0,0900	3,0531
7000	0,3741	0,1846	2,0266	0,1222	3,0615
8000	0,4894	0,2412	2,0286	0,1605	3,0491
9000	0,6186	0,3050	2,0286	0,2030	3,0475
10000	0,7637	0,3766	2,0281	0,2505	3,0483

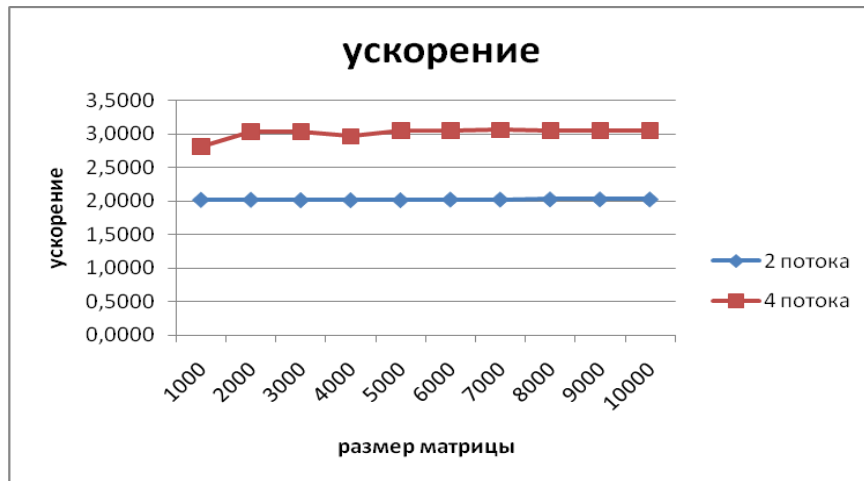


Рис. 6.7. Зависимость ускорения от количества исходных данных при выполнении параллельного алгоритма умножения матрицы на вектор, основанного на ленточном горизонтальном разбиении матрицы

Результаты вычислительных экспериментов приведены в табл. 6.3 и на рис. 6.7. Времена выполнения алгоритмов указаны в секундах. В табл. 6.4 и 6.5 и на рис. 6.8 и 6.9 представлены результаты сравнения времени выполнения T_p параллельного алгоритма умножения матрицы на вектор с использованием двух и четырех потоков, со временем T_p^* , полученным при помощи модели (6.9). Частота кэш-промахов, измеренная с помощью сис-

темы VPS, для двух потоков оказалась равной 0,0039, а для четырех потоков значение этой величины была оценена как 0,0143.

Таблица 6.4.

Сравнение экспериментального и теоретического времени выполнения параллельного алгоритма умножения матрицы на вектор, основанного на ленточном горизонтальном разбиении матрицы с использованием двух потоков

Размер матрицы	T_p	T_p^* (calc) (модель)	Модель 6.6 – оценка сверху		Модель 6.9 – уточненная оценка	
			T_p^* (mem)	T_p^*	T_p^* (mem)	T_p^*
1000	0,0038	0,0038	0,0131	0,0169	0,0001	0,0038
2000	0,0150	0,0151	0,0524	0,0675	0,0002	0,0153
3000	0,0341	0,0340	0,1179	0,1519	0,0005	0,0345
4000	0,0606	0,0605	0,2096	0,2701	0,0008	0,0613
5000	0,0947	0,0945	0,3275	0,4220	0,0013	0,0958
6000	0,1356	0,1361	0,4716	0,6077	0,0018	0,1379
7000	0,1846	0,1852	0,6420	0,8272	0,0025	0,1877
8000	0,2412	0,2419	0,8385	1,0804	0,0033	0,2452
9000	0,3050	0,3062	1,0612	1,3674	0,0041	0,3103
10000	0,3766	0,3780	1,3101	1,6881	0,0051	0,3831

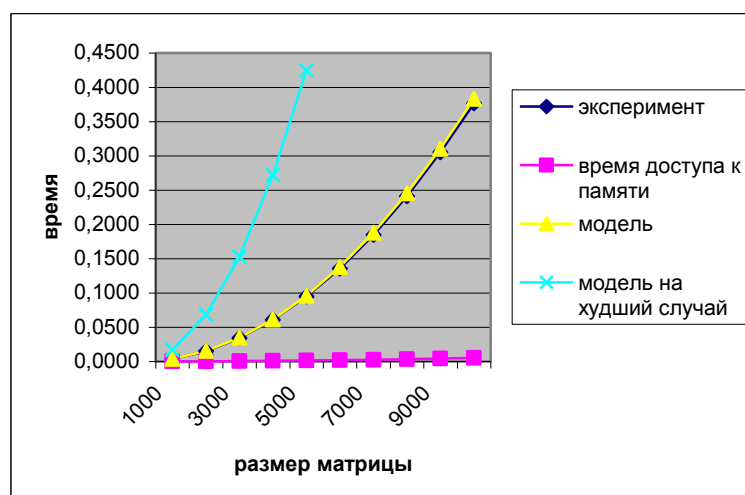


Рис. 6.8. График зависимости экспериментального и теоретического времени выполнения параллельного алгоритма от объема исходных данных при использовании двух потоков (ленточное разбиение матрицы по строкам)

Таблица 6.5.

Сравнение экспериментального и теоретического времени выполнения параллельного алгоритма умножения матрицы на вектор, основанного на ленточном горизонтальном разбиении матрицы с использованием четырех потоков

Размер матрицы	T_p	T_p^* (calc) (модель)	Модель 6.6 – оценка сверху		Модель 6.9 – уточненная оценка	
			T_p^* (mem)	T_p^*	T_p^* (mem)	T_p^*
1000	0,0027	0,0019	0,0131	0,0150	0,0002	0,0021
2000	0,0100	0,0076	0,0524	0,0600	0,0007	0,0083
3000	0,0227	0,0170	0,1179	0,1349	0,0017	0,0187
4000	0,0412	0,0302	0,2096	0,2399	0,0030	0,0332
5000	0,0626	0,0472	0,3275	0,3748	0,0047	0,0519
6000	0,0900	0,0680	0,4716	0,5397	0,0067	0,0748
7000	0,1222	0,0926	0,6420	0,7346	0,0092	0,1018
8000	0,1605	0,1210	0,8385	0,9594	0,0120	0,1329
9000	0,2030	0,1531	1,0612	1,2143	0,0152	0,1683
10000	0,2505	0,1890	1,3101	1,4991	0,0187	0,2077

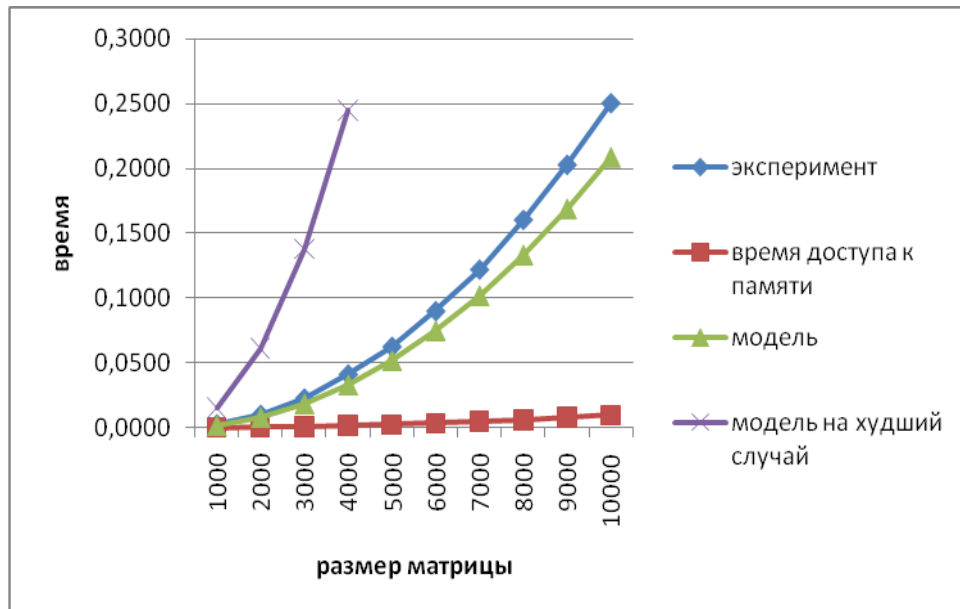


Рис. 6.9. График зависимости экспериментального и теоретического времени выполнения параллельного алгоритма от объема исходных данных при использовании четырех потоков (ленточное разбиение матрицы по строкам)

6.6. Умножение матрицы на вектор при разделении данных по столбцам

Рассмотрим теперь другой подход к параллельному умножению матрицы на вектор, основанный на разделении исходной матрицы на непрерывные вертикальные наборы (полосы) столбцов.

6.6.1. Определение подзадач и выделение информационных зависимостей

При таком способе разделения данных в качестве базовой подзадачи может быть выбрана операция умножения столбца матрицы A на один из элементов вектора b . Для организации вычислений в этом случае каждая базовая подзадача i , $0 \leq i < n$, должна иметь доступ к i -му столбцу матрицы A .

Каждая базовая задача i выполняет умножение своего столбца матрицы A на элемент b_i , в итоге в каждой подзадаче получается вектор $c'(i)$ промежуточных результатов. Для получения элементов результирующего вектора c необходимо просуммировать векторы $c'(i)$, которые были получены каждой подзадачей.

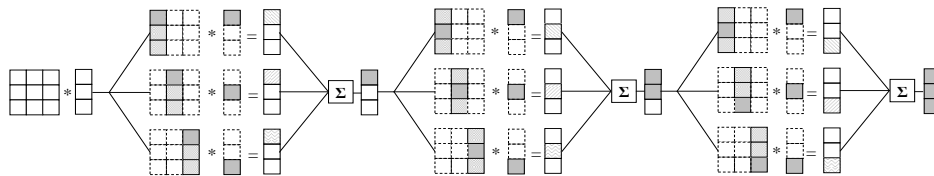


Рис. 6.10. Организация вычислений при выполнении параллельного алгоритма умножения матрицы на вектор с использованием разбиения матрицы по столбцам

Для того, чтобы реализовать такой подход, необходимо распараллелить внутренний цикл алгоритма умножения матрицы на вектор. В результате данного подхода параллельные области будут создаваться для вычисления каждого отдельного элемента результирующего вектора. Программа будет представляться в виде набора последовательных (*однопоточковых*) и параллельных (*многопоточковых*) участков. Подобный принцип организации параллелизма получил название «вилочного» (*fork-join*) или *пульсирующего параллелизма*. При выполнении многопоточкового участка каждый поток выполняет умножение одного элемента своего столбца исходной матрицы на один элемент вектора. Следующий за этим однопоточковый участок производит суммирование полученных результатов для того, чтобы вычислить элемент результирующего вектора (см. рис. 6.10). Таким образом, программа будет состоять из чередования n многопоточковых и n однопоточковых участков.

6.6.2. Масштабирование и распределение подзадач по вычислительным элементам

Выделенные базовые подзадачи характеризуются одинаковой вычислительной трудоемкостью и равным объемом передаваемых данных. В случае, когда количество столбцов матрицы превышает число процессоров, базовые подзадачи можно укрупнить, объединив в рамках одной подзадачи несколько соседних столбцов (в этом случае исходная матрица A разбивается на ряд вертикальных полос). При соблюдении равенства размера полос такой способ агрегации вычислений обеспечивает равномерность распределения вычислительной нагрузки по процессорам, составляющим многопроцессорную вычислительную систему.

6.6.3. Программная реализация

При реализации данного подхода к параллельному умножению матрицы на вектор каждый очередной многопоточковый участок параллельного алгоритма служит для вычисления одного элемента результирующего вектора. Очевидный способ реализовать подобный подход может состоять в следующем:

```
// Программа 6.4
// Функция параллельного матрично-векторного умножения
void ParallelResultCalculation(double* pMatrix,
    double* pVector, double* pResult, int Size) {
    int i, j;
    for (i=0; i<Size; i++) {
#pragma omp parallel for
        for (j=0; j<Size; j++)
            // !!! Внимание - см. приводимые далее пояснения
            pResult[i] += pMatrix[i*Size+j]*pVector[j];
    }
}
```

Выполним проверку правильности выполнения полученной программы – разработаем для этого функцию *TestResult*. Для контроля результатов выполним умножение исходной матрицы на вектор при помощи последовательного алгоритма, а затем сравним поэлементно векторы, полученные в результате выполнения последовательного и параллельного алгоритмов (следует отметить, что сравнение равенства вещественных чисел следует выполнять с некоторой явно указываемой точностью):

```
// Функция проверки матрично-векторного умножения
void TestResult(double* pMatrix, double* pVector,
```

```
double* pResult, int Size) {
double* pSerialResult = new double [Size];
double Accuracy = 1.0e-6;
int equal = 0;

// Последовательный алгоритм
SerialResultCalculation(pMatrix, pVector,
    pSerialResult, Size);

// Проверка результатов
for (int i=0; i<Size; i++)
    if (fabs(pResult[i]-pSerialResult[i])<=Accuracy)
        equal++;
if (equal<Size)
    printf("Результаты не совпадают \n");
else
    printf ("Результаты совпадают \n");

delete [] pSerialResult;
}
```

Результатом работы этой функции является печать диагностического сообщения. Данная функция позволяет контролировать правильность выполнения параллельного алгоритма в автоматическом режиме, независимо от сложности и объема исходных данных.

Результаты экспериментов показывают, что разработанная функция параллельного умножения матрицы на вектор, реализующая разделение данных по столбцам, не всегда дает верный результат. Причина неправильной работы кроется в неправильном использовании переменных, являющихся общими для нескольких потоков. В разработанном варианте программы такими общими переменными являются элементы вектора *pResult*. В ходе вычислений несколько разных потоков могут, например, попытаться одновременно записать новые значения в одну и ту же общую переменную, что приведет к получению ошибочных результатов. Для исключения взаимовлияния потоки должны использовать общие данные с соблюдением правил взаимного исключения, которые обеспечивали бы использование общих переменных в каждый момент времени только одним потоком (а потоки, пытающиеся получить доступ к «занятым» общим данным, должны блокироваться).

Реализация взаимного исключения доступа может состоять в использовании предусмотренных в OpenMP специальных переменных синхронизации — *замков*:

```
// Программа 6.5
// Функция параллельного матрично-векторного умножения
void ParallelResultCalculation(double* pMatrix,
    double* pVector, double* pResult, int Size) {
    int i, j;
    omp_lock_t MyLock;
    omp_init_lock(&MyLock);

    for (i=0; i<Size; i++) {
#pragma omp parallel for
        for (j=0; j<Size; j++) {
            omp_set_lock(&MyLock);
            pResult[i] += pMatrix[i*Size+j]*pVector[j];
            omp_unset_lock(&MyLock);
        }
    }
    omp_destroy_lock(&MyLock);
}
```

Функция *omp_set_lock* выполняется для потока только в том случае, если другие потоки не выполнили вызов этой функции для той же самой переменной; в противном случае поток блокируется (переменная-замок пропускает только один поток). Продолжение блокированных потоков осуществляется только после выполнения для замка функции *omp_unset_lock* (при наличии нескольких блокированных потоков после снятия замка выполнение разрешается только для одного потока и замок снова «закрывается»).

После внесения этих изменений в код функции, выполняющей умножение матрицы на вектор, результирующий вектор совпадает с вектором, полученным при помощи последовательного алгоритма. Но реализация алгоритма регулирования доступа к элементу вектора *pResult[i]* привела к тому, что потоки многопоточкового участка могут выполняться только последовательно, т. к. один поток не может получить право на изменение разделяемой переменной до тех пор, пока такое изменение выполняет другой поток.

Для одновременного решения проблемы разделения и защиты общих переменных в OpenMP реализован стандартный механизм эффективного выполнения операции редукции при распараллеливании циклов. Данный механизм состоит в использовании директивы *reduction*. Если при распараллеливании циклов указывается директива *reduction*, компилятор автоматически создает локальные копии переменной редукции для каждого потока параллельной программы, а после вы-

полнения цикла собирает значения локальных копий в разделяемую переменную с использованием операции, указанной как аргумент директивы редукции.

Представим новый вариант параллельной программы умножения матрицы на вектор с использованием алгоритма разбиения матрицы по столбцам. Ниже приведем только код основной функции, выполняющей параллельный алгоритм умножения матрицы на вектор.

```
// Программа 6.6
// Функция параллельного матрично-векторного умножения
void ParallelResultCalculation(double* pMatrix,
    double* pVector, double* pResult, int Size) {
    int i, j;
    double IterGlobalSum = 0;

    for (i=0; i<Size; i++) {
        IterGlobalSum = 0;
        #pragma omp parallel for reduction(+:IterGlobalSum)
        for (j=0; j<Size; j++)
            IterGlobalSum += pMatrix[i*Size+j]*pVector[j];
        pResult[i] = IterGlobalSum;
    }
}
```

6.6.4. Анализ эффективности

При анализе эффективности данного параллельного алгоритма будем полагаться на результаты, полученные в п. 6.5.4. Очевидно, что ни время выполнения одной вычислительной операции, ни эффективная скорость доступа к оперативной памяти не зависят от того, каким образом распределяются вычисления между потоками параллельной программы.

Прежде всего следует заметить, что при выполнении параллельного алгоритма умножения матрицы на вектор, основанного на вертикальном разделении матрицы, выполняется больше арифметических операций. Действительно, на каждой итерации внешнего цикла после вычисления каждым потоком своей частичной суммы необходимо выполнить редукцию полученных результатов. Сложность выполнения операции редукции — $\log_2 p$. После выполнения редукции данных необходимо запомнить результат в очередной элемент результирующего вектора. Таким образом, время вычислений для данного алгоритма определяется формулой:

$$T_{calc} = \left(\frac{n \cdot 2n - 1}{p} + n \cdot \log_2 p + n \right) \cdot \tau.$$

Как видно из представленного программного кода, на каждой итерации внешнего цикла алгоритма умножения матрицы на вектор организуется параллельная секция для вычисления каждого элемента результирующего вектора. Для организации и синхронизированного закрытия параллельных секций, а также для выполнения редукции, вызываются специализированные функции библиотеки OpenMP. С одной стороны, для выполнения каждой такой функции необходимо определенное время. С другой стороны, в процессе выполнения эти функции используют данные, которые должны быть загружены в кэш процессора. Так как объем кэша ограничен, это ведет к вытеснению и последующей повторной загрузке элементов матрицы и вектора. На выполнение таких подкачек данных тоже тратится некоторое время.

Для анализа программного кода и подтверждения выдвинутых предположений воспользуемся специализированной программной системой, предназначенной для профилирования и оптимизации программ, Intel VTune Performance Analyzer (см., например, [107]). Данная программная система позволяет находить в программе критический (вычислительно-трудоемкий) код, измерять различные характеристики эффективности кода и наблюдать за процессом выполнения программы в динамике.

Проанализируем при помощи инструмента Call Graph (*Граф вызова функций*) профилировщика Intel VTune Performance Analyzer приложение, выполняющее параллельный алгоритм умножения матрицы на вектор, основанный на разделении матрицы на горизонтальные полосы, и сконцентрируем свое внимание на анализе выполнения функции *ParallelResultCalculation*. На рис. 6.11 представлен результат анализа. Видно, что собственное (без учета времени выполнения вызываемых функций) время и полное время выполнения функции совпадают. Это означает, что, несмотря на то, что функция *ParallelResultCalculation* вызывает другие функции (в данном случае – функции библиотеки времени исполнения OpenMP), практически все время тратится непосредственно на вычисления.

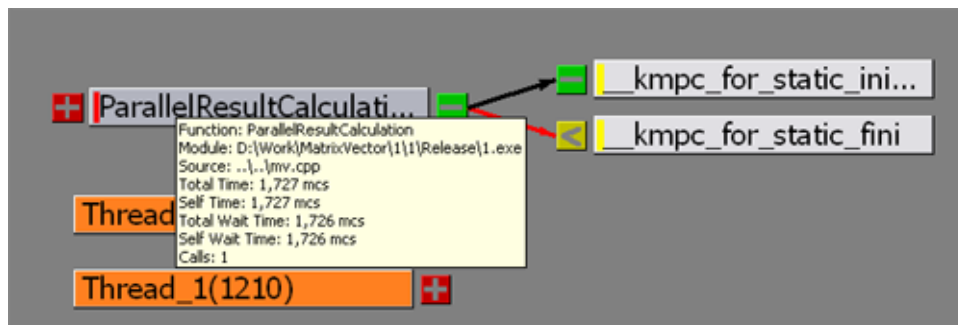


Рис. 6.11. Общее и собственное время выполнения функции *ParallelResultCalculation* в случае разделения матрицы на горизонтальные полосы

Теперь перейдем к анализу параллельного алгоритма, основанного на разделении матрицы по столбцам. В данном случае собственное время выполнения функции *ParallelResultCalculation* существенно меньше полного времени ее выполнения. Большая доля времени тратится на выполнение вложенных функций, которыми являются функции, необходимые для организации и закрытия параллельных секций, а также для выполнения редукции (рис. 6.12).

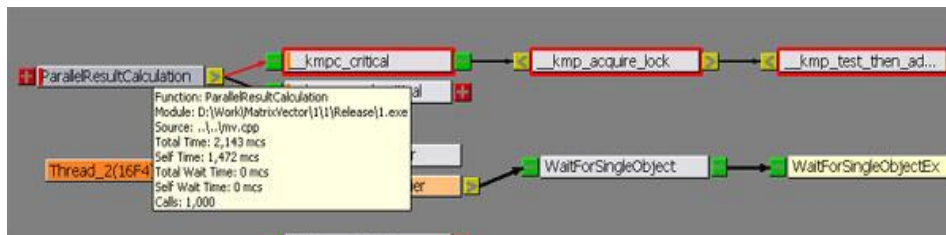


Рис. 6.12. Общее и собственное время выполнения функции *ParallelResultCalculation* в случае разделения матрицы на вертикальные полосы

Далее воспользуемся инструментом Sampling³⁾ профилировщика Intel VTune Performance Analyzer, обеспечивающим автоматизированный процесс регистрации различных событий, возникающих в процессоре во время исполнения профилируемой программы. Создадим проект, в котором вызовем последовательно функцию, выполняющую параллельный алгоритм умножения матрицы на вектор, основанный на горизонтальном разделении

³⁾ В научно-технической литературе практически не встречается общепринятого именования инструмента Sampling на русском языке и для его обозначения часто используют транслитерацию вида *Сэмплирование*.

матрицы (*ParallelResultCalculation1*), и функцию, выполняющую параллельный алгоритм, основанный на вертикальном разделении матрицы (*ParallelResultCalculation2*). Измерим число возникновений события выделения новой кэш-строки первого уровня (в Intel VTune Performance Analyzer данное событие называется L1 Cache Line Allocation). Результаты профилирования приложения представлены на рис. 6.13. Эксперименты проводились для матрицы размером 1000×1000 элементов.

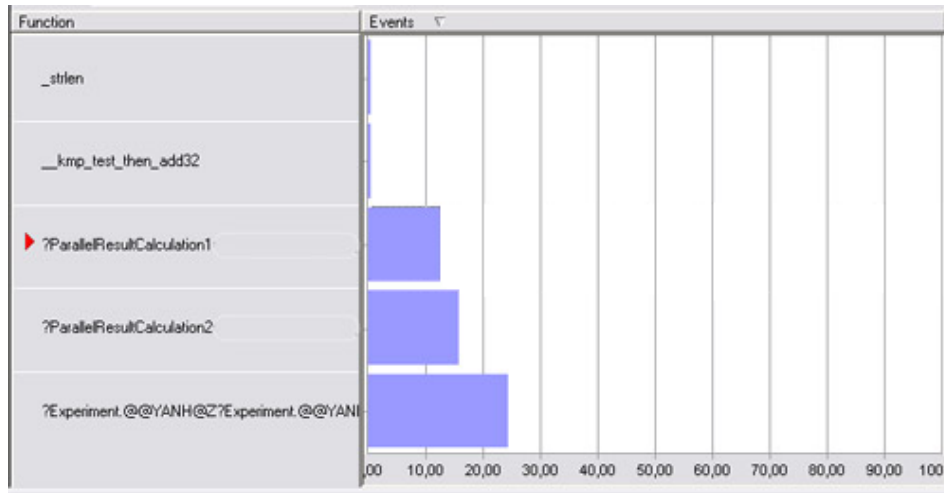


Рис. 6.13. Соотношение количества событий выделения кэш-линий первого уровня для двух параллельных алгоритмов умножения матрицы на вектор

При выполнении параллельного алгоритма умножения матрицы на вектор, основанного на вертикальном разделении, выделяется больше (в среднем на 25%) строк кэша первого уровня.

Теперь оценим накладные расходы на организацию параллельности на каждой итерации алгоритма, основанного на вертикальном разделении данных. Для этого подготовим еще один вариант параллельного алгоритма умножения матрицы на вектор при разделении матрицы по столбцам. Образум в этом варианте программы потоки, которые самостоятельно определяют (используя идентификатор потока) блоки элементов матрицы и векторы для обработки. Далее потоки параллельно выполняют умножение своих столбцов матрицы на элементы вектора, результат умножения сохраняется в общем массиве *AllResults*. Количество элементов в этом массиве равно $Size \times ThreadNum$, где *ThreadNum* – общее число потоков. Потоки осуществляют запись частичных результатов в непересекающиеся сегменты массива *AllResults*: нулевой поток осуществляет запись в элементы с 0 по $Size-1$, первый поток – в элементы с $Size$ по $2 \cdot Size-1$, и т. д.

После выполнения умножения элементы вектора *AllResults* необходимо просуммировать для получения элементов результирующего вектора.

```
void ParallelResultCalculation(double* pMatrix,
    double* pVector, double* pResult,
    double* AllResults, int Size) {
    int ThreadNum;
#pragma omp parallel shared (ThreadNum)
    {
        ThreadNum = omp_get_num_threads();
        int ThreadID = omp_get_thread_num();
        int BlockSize = Size/ThreadNum;
        double IterResult;
        for (int i=0; i<Size; i++) {
            IterResult = 0;
            for (int j=0; j<BlockSize; j++)
                IterResult += pMatrix[i*Size+j+
                    ThreadID*BlockSize] *
                    pVector[j+ThreadID*BlockSize];
            AllResults[Size*ThreadID+i] = IterResult;
        }
    }
    for (int i=0; i<Size; i++)
        for (int j=0; j<ThreadNum; j++)
            pResult[i] += AllResults[j*Size+i];
}
```

Как видно из представленного программного кода, при выполнении данного алгоритма параллельная секция создается только один раз, и, следовательно, время, необходимое для выполнения функций библиотеки OpenMP, отвечающих за организацию и закрытие параллельных секций, пренебрежимо мало. Чтобы оценить накладные расходы δ на организацию параллельности на каждой итерации алгоритма, необходимо из времени выполнения исходного параллельного алгоритма умножения матрицы на вектор, основанного на вертикальном разделении данных, вычесть время выполнения вновь разработанного параллельного алгоритма и поделить полученную разницу на количество создаваемых параллельных секций. Эксперименты показывают, что величина δ равна 0.25 мкс. Таким образом, время выполнения рассматриваемого параллельного алгоритма умножения матрицы на вектор может быть вычислено для наихудшего случая (без использования кэш памяти) по формуле:

$$T_p = \left(\frac{n \cdot 2n-1}{p} + n \cdot \log_2 p + n \right) \cdot \tau + n^2 \cdot \left(\alpha + \frac{64}{\beta} \right) + n \cdot \delta. \quad (6.10)$$

При использовании кэш памяти с частотой кэш-промахов γ соотношение (6.10) принимает вид:

$$T_p = \left(\frac{n \cdot 2n-1}{p} + n \cdot \log_2 p + n \right) \cdot \tau + \gamma \cdot n^2 \cdot \left(\alpha + \frac{64}{\beta} \right) + n \cdot \delta. \quad (6.11)$$

6.6.5. Результаты вычислительных экспериментов

Вычислительные эксперименты для оценки эффективности параллельного алгоритма умножения матрицы на вектор при разбиении данных по столбцам проводились при условиях, указанных в п. 6.5.5. Результаты вычислительных экспериментов приведены в табл. 6.6.

Таблица 6.6.

Результаты вычислительных экспериментов по исследованию параллельного алгоритма умножения матрицы на вектор, основанного на разбиении матрицы по столбцам

Размер матрицы	Последовательный алгоритм	Параллельный алгоритм			
		2 потока		4 потока	
		Время	Ускорение	время	Ускорение
1000	0,0076	0,0067	1,1335	0,0064	1,1929
2000	0,0303	0,0219	1,3825	0,0165	1,8349
3000	0,0688	0,0452	1,5218	0,0316	2,1748
4000	0,1222	0,0769	1,5897	0,0503	2,4292
5000	0,1909	0,1160	1,6465	0,0712	2,6812
6000	0,2748	0,1624	1,6918	0,0979	2,8081
7000	0,3741	0,2186	1,7116	0,1340	2,7918
8000	0,4894	0,2819	1,7359	0,1641	2,9818
9000	0,6186	0,3539	1,7479	0,2036	3,0380
10000	0,7637	0,4330	1,7639	0,2479	3,0812

В табл. 6.7, 6.8 и на рис. 6.15, 6.16 представлены результаты сравнения времени выполнения T_p параллельного алгоритма умножения матрицы на вектор с использованием двух и четырех потоков со временем T_p^* , полученным при помощи модели (6.11). Частота кэш-промахов, измеренная с помощью системы VPS, для двух потоков оказалась равной 0,0039, а для четырех потоков значение этой величины была оценена как 0,0377.

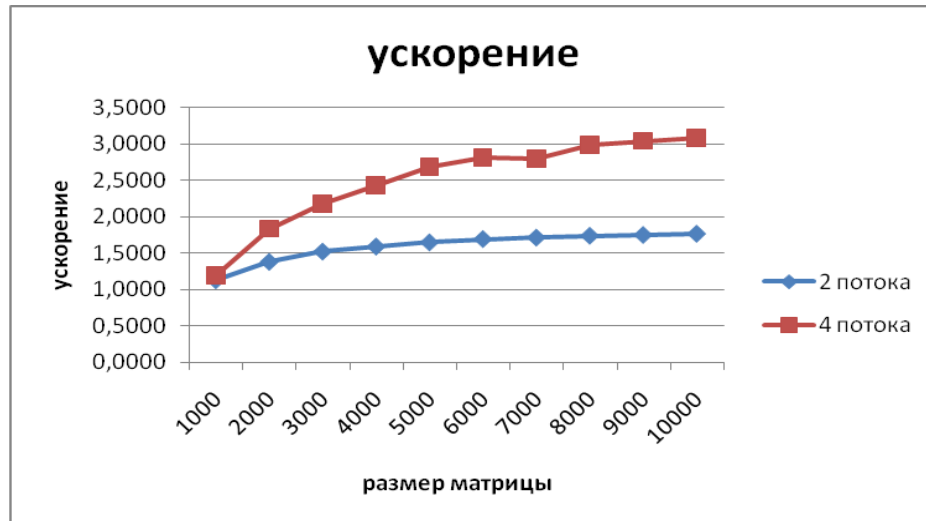


Рис. 6.14. Зависимость ускорения размера матрицы при выполнении параллельного алгоритма умножения матрицы на вектор (ленточное разбиение матрицы по столбцам)

Таблица 6.7.

Сравнение экспериментального и теоретического времени выполнения параллельного алгоритма умножения матрицы на вектор, основанного на ленточном разбиении матрицы по столбцам с использованием двух потоков

Размер матрицы	T_p	$T_p^* (calc)$ (модель)	Модель 6.10 – оценка сверху		Модель 6.11 – уточненная оценка	
			$T_p^* (mem)$	T_p^*	$T_p^* (mem)$	T_p^*
1000	0,0067	0,0040	0,0131	0,0171	0,0001	0,0041
2000	0,0219	0,0156	0,0524	0,0680	0,0002	0,0158
3000	0,0452	0,0348	0,1179	0,1527	0,0005	0,0352
4000	0,0769	0,0615	0,2096	0,2711	0,0008	0,0623
5000	0,1160	0,0958	0,3275	0,4233	0,0013	0,0971
6000	0,1624	0,1376	0,4716	0,6093	0,0018	0,1395
7000	0,2186	0,1870	0,6420	0,8290	0,0025	0,1895
8000	0,2819	0,2440	0,8385	1,0825	0,0033	0,2472
9000	0,3539	0,3085	1,0612	1,3697	0,0041	0,3126
10000	0,4330	0,3806	1,3101	1,6907	0,0051	0,3857

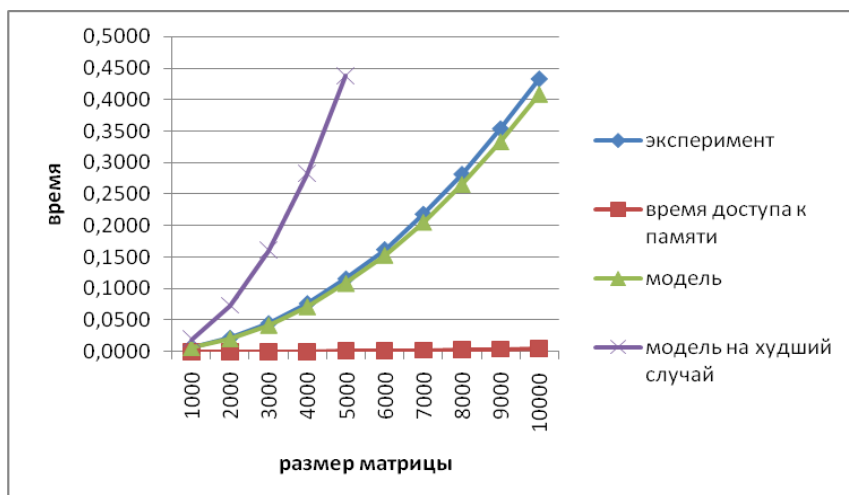


Рис. 6.15. График зависимости экспериментального и теоретического времен выполнения параллельного алгоритма от объема исходных данных при использовании двух потоков (ленточное разбиение матрицы по столбцам)

Таблица 6.8.
Сравнение экспериментального и теоретического времен выполнения параллельного алгоритма умножения матрицы на вектор, основанного на ленточном разбиении матрицы по столбцам с использованием четырех потоков

Размер матрицы	T_p	$T_p^* (calc)$ (модель)	Модель 6.10 – оценка сверху		Модель 6.11 – уточненная оценка	
			$T_p^* (mem)$	T_p^*	$T_p^* (mem)$	T_p^*
1000	0,0064	0,0022	0,0131	0,0153	0,0005	0,0026
2000	0,0165	0,0081	0,0524	0,0605	0,0020	0,0101
3000	0,0316	0,0178	0,1179	0,1357	0,0044	0,0222
4000	0,0503	0,0313	0,2096	0,2409	0,0079	0,0392
5000	0,0712	0,0486	0,3275	0,3761	0,0123	0,0609
6000	0,0979	0,0696	0,4716	0,5413	0,0178	0,0874
7000	0,1340	0,0944	0,6420	0,7364	0,0242	0,1186
8000	0,1641	0,1230	0,8385	0,9615	0,0316	0,1547
9000	0,2036	0,1554	1,0612	1,2166	0,0400	0,1954
10000	0,2479	0,1916	1,3101	1,5017	0,0494	0,2410

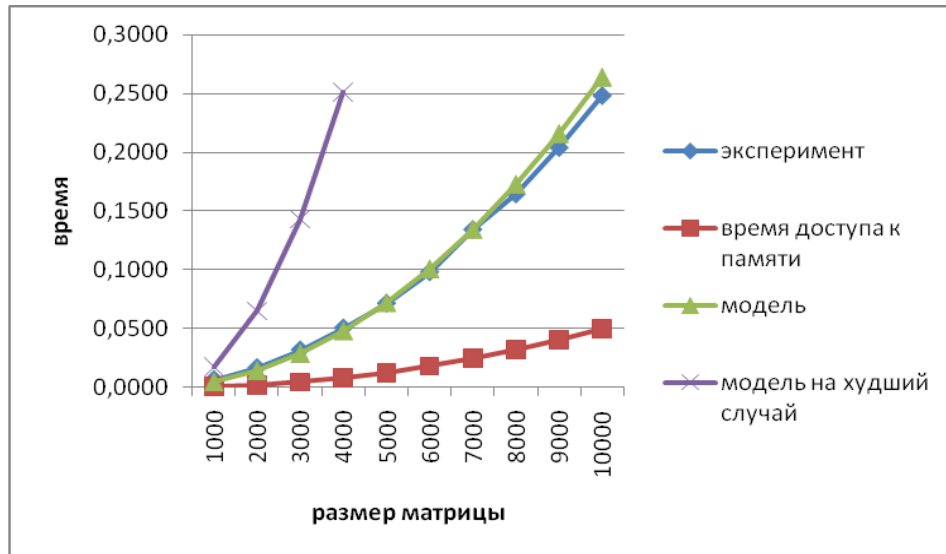


Рис. 6.16. График зависимости экспериментального и теоретического времени выполнения параллельного алгоритма от объема исходных данных при использовании четырех потоков (ленточное разбиение матрицы по столбцам)

6.7. Умножение матрицы на вектор при блочном разделении данных

Рассмотрим теперь параллельный алгоритм умножения матрицы на вектор, который основан на ином способе разделения данных – на разбиении матрицы на прямоугольные фрагменты (*блоки*).

6.7.1. Определение подзадач

Блочная схема разбиения матриц подробно рассмотрена в разделе 6.2. При таком способе разделения данных исходная матрица A представляется в виде набора прямоугольных блоков:

$$A = \begin{pmatrix} A_{00} & A_{02} & \dots A_{0q-1} \\ & \dots & \\ A_{s-11} & A_{s-12} & \dots A_{s-1q-1} \end{pmatrix},$$

где A_{ij} , $0 \leq i < s$, $0 \leq j < q$, есть блок матрицы:

$$A_{ij} = \begin{pmatrix} a_{i_0j_0} & a_{i_0j_1} & \dots a_{i_0j_{l-1}} \\ & \dots & \\ a_{i_{k-1}j_0} & a_{i_{k-1}j_1} & a_{i_{k-1}j_{l-1}} \end{pmatrix}, \quad i_v = ik + v, 0 \leq v < k, k = m/s, \\ j_u = jl + u, 0 \leq u \leq l, l = n/q$$

(здесь, как и ранее, предполагается, что $p = s \cdot q$, количество строк матрицы является кратным s , а количество столбцов – кратным q , то есть $m = k \cdot s$ и $n = l \cdot q$).

При использовании блочного представления матрицы A базовые подзадачи целесообразно определить на основе вычислений, выполняемых над матричными блоками. Для нумерации подзадач могут использоваться индексы располагаемых в подзадачах блоков матрицы A , т. е. подзадача (i, j) производит вычисления над блоком A_{ij} . Помимо блока матрицы A каждая подзадача должна иметь доступ к блоку вектора b . При этом для блоков одной и той же подзадачи должны соблюдаться определенные правила соответствия – операция умножения блока матрицы A_{ij} может быть выполнена только, если блок вектора $b'(i, j)$ имеет вид

$$b'(i, j) = (b'_0(i, j), \dots, b'_{l-1}(i, j)), \text{ где} \\ b'_u(i, j) = b_{j_u}, j_u = jl + u, 0 \leq u < l, l = n/q.$$

6.7.2. Выделение информационных зависимостей

Рассмотрим общую схему параллельных вычислений для операции умножения матрицы на вектор при блочном разделении исходных данных. После перемножения блоков матрицы A и вектора b каждая подзадача (i, j) будет содержать вектор частичных результатов $c'(i, j)$, определяемый в соответствии с выражениями:

$$c'_v(i, j) = \sum_{u=0}^{l-1} a_{i_vj_u} b_{j_u}, i_v = ik + v, 0 \leq v < k, k = m/s, \\ j_u = jl + u, 0 \leq u \leq l, l = n/q.$$

Как можно видеть из приведенных выражений, каждый поток вычисляет блок вектора частичных результатов исходной задачи умножения матрицы на вектор. Полный результат – вектор c – можно определить суммированием элементов блока вектора частичных результатов с одинаковыми индексами по всем подзадачам, относящимся к одним и тем же строкам решетки потоков, т. е.

$$c_v(i) = \sum_{j=0}^{q-1} c'_v(i, j), 0 \leq i < s, 0 \leq v < k, k = m / s.$$

Общая схема выполняемых вычислений для умножения матрицы на вектор при блочном разделении данных показана на рис. 6.17.

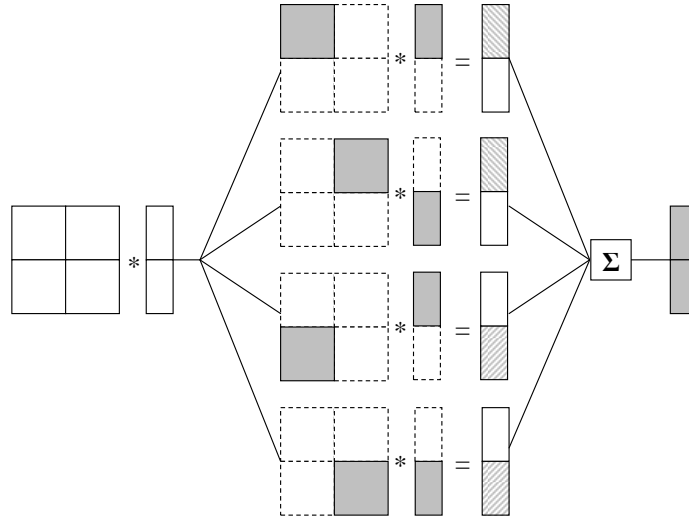


Рис. 6.17. Общая схема параллельного алгоритма умножения матрицы на вектор при блочном разделении данных

При дальнейшем изложении материала для снижения сложности и упрощения получаемых соотношений будем полагать, что число блоков в матрице A совпадает по горизонтали и по вертикали, т. е. $s = q$. Для эффективного выполнения параллельного алгоритма умножения матрицы на вектор, основанного на блочном разделении данных, целесообразно выделить число параллельных потоков совпадающим с количеством блоков матрицы A , т. е. такое количество потоков, которое является полным квадратом q^2 . Данное требование может привести к тому, что число вычислительных элементов и количество потоков может различаться – для обозначения числа потоков в дальнейшем будем использовать переменную π , т. е. $\pi = q^2$. Дополнительно можно отметить заранее, что для эффективного выполнения вычислений количество потоков π должно быть, по крайней мере, кратным числу вычислительных элементов p .

На рис. 6.17 приведена схема информационного взаимодействия подзадач в случае, когда для выполнения алгоритма создано 4 потока.

Рассмотрев представленную схему параллельных вычислений, можно сделать вывод о том, что информационная зависимость базовых подзадач

проявляется только на этапе суммирования результатов перемножения блоков матрицы A и блоков вектора b .

6.7.3. Масштабирование и распределение подзадач по вычислительным элементам

При сделанных предположениях общее количество базовых подзадач совпадает с числом выделенных потоков. Так, если определить число потоков $\pi = q \cdot q$, то размер блоков матрицы A определяется соотношениями:

$$k=m/q, \quad l=n/q,$$

где k и l есть количество строк и столбцов в блоках матрицы A . Такой способ определения размера блоков приводит к тому, что объем вычислений в каждой подзадаче является равным и тем самым достигается полная балансировка вычислительной нагрузки между потоками.

6.7.4. Программная реализация

1. Первый вариант. Как уже отмечалось выше, для эффективного выполнения алгоритма умножения матрицы на вектор необходимо, чтобы количество параллельных потоков являлось полным квадратом. В приведенном ниже примере используется 4 потока (значение переменной *GridThreadsNum* устанавливается равным 4 – значение данной переменной должно переустанавливаться при изменении необходимого числа потоков). Определение числа потоков «вручную» до начала выполнения программы гарантирует корректную работу приложения на вычислительной системе с любым количеством доступных вычислительных элементов. Однако, если в вычислительной системе больше вычислительных элементов, чем определено параллельных потоков, эффективное использование ресурсов не может быть достигнуто.

Для определения количества параллельных потоков используется функция *omp_set_num_threads* библиотеки OpenMP (функция должна быть вызвана в последовательном участке программы):

```
int omp_set_num_threads (int NumThreads);
```

Чтобы корректно определять блоки данных, над которыми поток должен выполнять вычисления, в программе необходимо иметь уникальный идентификатор потока. В качестве такого идентификатора может выступать номер потока. В библиотеке OpenMP существует функция для определения номера потока:

```
int omp_get_thread_num (void);
```

В представленном варианте алгоритма номер потока сохраняется в переменной *ThreadID*, которая при объявлении параллельной секции определена как *private*, т. е. для этой переменной в каждом потоке создается локальная копия. Положения блоков матрицы и вектора, которые должны обрабатываться потоком, определяются в зависимости от значения переменной *ThreadID*.

Для накопления результатов умножения блоков матрицы и вектора в каждом потоке используется локальная область памяти *pThreadResult*. Для того, чтобы получить элемент результирующего вектора, необходимо просуммировать соответствующие элементы векторов *pThreadResult* со всех потоков. Для обеспечения контроля доступа к разделяемому ресурсу *pResult* со всех потоков параллельной программы используется механизм *критических секций*. Для объявления критической секции служит директива:

```
#pragma omp critical  
    structured block
```

Блок операций, следующий за объявлением критической секции, в каждый момент времени может выполняться только одним потоком. Подобная организация программы гарантирует единственность доступа потоков для изменения разделяемых данных внутри критической секции.

Представим возможный вариант параллельной программы умножения матрицы на вектор с использованием алгоритма разбиения матрицы по блокам. Ниже приведен только код основной функции, выполняющей параллельный алгоритм умножения матрицы на вектор.

```
// Программа 6.7  
// Функция параллельного матрично-векторного умножения  
void ResultCalculation(double* pMatrix,  
    double* pVector, double* pResult, int Size) {  
    int ThreadID;  
    int GridThreadsNum = 4;  
    int GridSize = int (sqrt(double(GridThreadsNum)));  
    // Предполагается, что размер матрицы кратен  
    // размеру сетки потоков  
    int BlockSize = Size/GridSize;  
  
    omp_set_num_threads(GridThreadsNum);  
  
#pragma omp parallel private (ThreadID)  
    {  
        ThreadID = omp_get_thread_num();  
        double * pThreadResult = new double [Size];  
        for (int i=0; i<Size; i++) {
```

```

    pThreadResult[i] = 0;
}
int i_start = (int(ThreadID/GridSize))*BlockSize;
int j_start = (ThreadID%GridSize)*BlockSize;
double IterResult;
for (int i=0; i< BlockSize; i++) {
    IterResult = 0;
    for (int j=0; j < BlockSize; j++)
        IterResult +=
            pMatrix[(i+i_start)*Size+(j+j_start)]*
            pVector[j+i_start];
    pThreadResult[i+i_start] = IterResult;
}
#pragma omp critical
for (int i=0; i<Size; i++) {
    pResult[i] += pThreadResult[i];
}
delete [] pThreadResult;
} // pragma omp parallel
}

```

Следует отметить, что для запоминания результатов потоков в *pThreadResult* и их объединения в массиве *pResult* используется несколько «упрощенная» схема реализации с целью получения более простого варианта программного кода.

2. Второй вариант. При реализации первого варианта блочного параллельного алгоритма умножения матрицы на вектор необходимо «вручную» определить блоки данных и блоки вектора, над которыми каждый поток выполняет вычисления. Однако стандарт OpenMP предусматривает возможность организации *вложенных параллельных секций*. Это значит, что внутри одной параллельной секции можно объявить вторую параллельную секцию, которая разделит каждый из потоков внешней параллельной секции на несколько вложенных потоков.

Продemonстрируем описанный подход в случае, когда при объявлении каждой новой параллельной секции поток выполнения разделяется на два. При перемножении матрицы на вектор для внешнего цикла объявим внешнюю параллельную секцию. Потоки этой секции разделяют между собой строки матрицы – каждый поток в своих вычислениях использует горизонтальную полосу матрицы. Далее, используя механизм вложенного параллелизма, объявим внутреннюю параллельную секцию: потоки этой параллельной секции делят между собой столбцы полосы матрицы. Таким образом, будет реализовано блочное разделение данных. Схема выполнения данного параллельного алгоритма представлена на рис. 6.18.

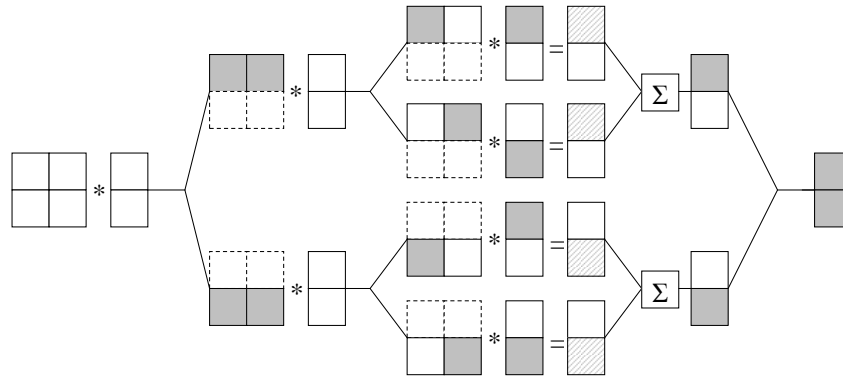


Рис. 6.18. Организация вычислений при выполнении параллельного алгоритма умножения матрицы на вектор с использованием разбиения матрицы на блоки и вложенного параллелизма

Использование механизма OpenMP для организации вложенного параллелизма позволяет значительно проще реализовать параллельный алгоритм умножения матрицы на вектор, основанный на блочном разделении матриц. Однако следует отметить, что на данный момент не все компиляторы, реализующие стандарт OpenMP, поддерживают вложенный параллелизм. Для компиляции представленного ниже кода использовался компилятор Intel C++ Compiler 10.0 for Windows, поддерживающий вложенный параллелизм.

Для того, чтобы включить поддержку вложенного параллелизма, необходимо вызвать функцию `omp_set_nested` библиотеки OpenMP (функция должна быть вызвана из последовательного участка программы):

```
void omp_set_nested(int nested);
```

Для разделения матрицы на полосы необходимо разделить внешний цикл алгоритма умножения матрицы на вектор при помощи директивы `omp parallel for`. Для дальнейшего разделения полос матрицы на блоки – разделить внутренний цикл между потоками при помощи той же директивы. Поскольку для вычисления каждого элемента результирующего вектора создается параллельная секция, внутри которой каждый поток вычисляет частичную сумму для этого элемента, необходимо применить механизм редукции, подробно описанный в разделе 6.6.

// Программа 6.8

```
// Функция параллельного матрично-векторного умножения
void ResultCalculation(double* pMatrix,
    double* pVector, double* pResult,
    int Size) {
    int NestedThreadsNum = 2;
```

```

omp_set_num_threads(NestedThreadsNum);
omp_set_nested(true);
#pragma omp parallel for
for (int i=0; i<Size; i++) {
    double ThreadResult = 0;
#pragma omp parallel for reduction(+:ThreadResult)
    for (int j=0; j<Size; j++)
        ThreadResult += pMatrix[i*Size+j]*pVector[j];
    pResult[i] = ThreadResult;
}
}

```

Отметим, что в приведенной программе для задания количество потоков, создаваемых на каждом уровне вложенности параллельных областей, используется переменная *NestedThreadsNum* (в данном варианте программы ее значение устанавливается равным 2 – данное значение должно переставляться при изменении необходимого числа потоков).

6.7.5. Анализ эффективности

1. Первый вариант реализации. При анализе эффективности первой реализации параллельного алгоритма умножения матрицы на вектор, основанного на блочном разделении матрицы, обратим внимание на дополнительные вычислительные операции: после того, как каждый поток выполнил умножение блока матрицы на блок вектора, необходимо сложить вектора частичных результатов для получения результирующего вектора. Если для выполнения параллельного алгоритма использовалось π потоков, то, с учетом использованной при реализации алгоритма схемы редукции данных, количество дополнительных операций равно $n \cdot \pi$. Таким образом, время вычислений можно определить по формуле:

$$T_{calc} = \left(\frac{n \cdot 2n - 1}{p} + n \cdot \pi \right) \cdot \tau.$$

Для оценки полного времени выполнения параллельного алгоритма можно использовать все положения, использованные при выводе необходимых аналитических соотношений для параллельного алгоритма при ленточном горизонтальном разделении данных (см. п. 6.5.4). Как результат, оценка сложности параллельных вычислений для наихудшего случая может быть представлена в следующем виде:

$$T_p = \left(\frac{n \cdot 2n - 1}{p} + n \cdot \pi \right) \cdot \tau + n^2 \cdot \left(\alpha + \frac{64}{\beta} \right). \quad (6.12)$$

При учете частоты кэш промахов соотношение (6.12) имеет вид:

$$T_p = \left(\frac{n \cdot 2n-1}{p} + n \cdot \pi \right) \cdot \tau + \gamma \cdot n^2 \cdot \left(\alpha + \frac{64}{\beta} \right). \quad (6.13)$$

Инструмент Call Graph профилировщика Intel VTune Performance Analyzer показывает, что затраты на организацию параллелизма не превосходят 5%.

2. Второй вариант реализации. Выполнив анализ второй вариант программной реализации параллельного алгоритма умножения матрицы на вектор, можно выделить дополнительные вычислительные операции: после вычисления частичного результата каждым параллельным потоком «внутренней» параллельной секции необходимо выполнить редукцию и записать полученный результат в соответствующий элемент результирующего вектора. Напомним, что для выполнения задачи используется π параллельных потоков, $\pi = q \cdot q$, и их число может отличаться от количества имеющихся вычислительных элементов (процессоров или ядер). С учетом данного обстоятельства можно определить, что время выполнения вычислений ограничено сверху величиной:

$$T_{calc} = \left(\frac{n \cdot 2n-1}{p} + n \cdot q \right) \cdot \tau.$$

В случае же, когда количества вычислительных элементов совпадает с числом потоков, т. е. $p = \pi$, оценка времени вычислений может быть получена более точно:

$$T_{calc} = \left(\frac{n \cdot 2n-1}{p} + \frac{n}{q} \cdot \log_2 q + \frac{n}{q} \right) \cdot \tau$$

При выполнении вычислений «внутренние» параллельные секции создаются и закрываются много раз; кроме того, дополнительное время тратится на синхронизацию и выполнение операции редукции. Для оценки доли накладных расходов на обеспечение параллелизма снова воспользуемся профилировщиком Intel VTune Performance Analyzer (рис. 6.19).

Накладные расходы на организацию и закрытие параллельных секций были измерены при анализе параллельного алгоритма, основанного на вертикальном разделении данных (см. раздел 6.6). В данном случае каждый поток создает параллельные секции n/q раз. Таким образом, время выполнения параллельного алгоритма может быть вычислено по формуле:

$$T_p = \left(\frac{n \cdot 2n-1}{p} + \frac{n}{q} \cdot \log_2 q + \frac{n}{q} \right) \cdot \tau + n^2 \cdot \left(\alpha + \frac{64}{\beta} \right) + \frac{n}{q} \cdot \delta \quad (6.14)$$

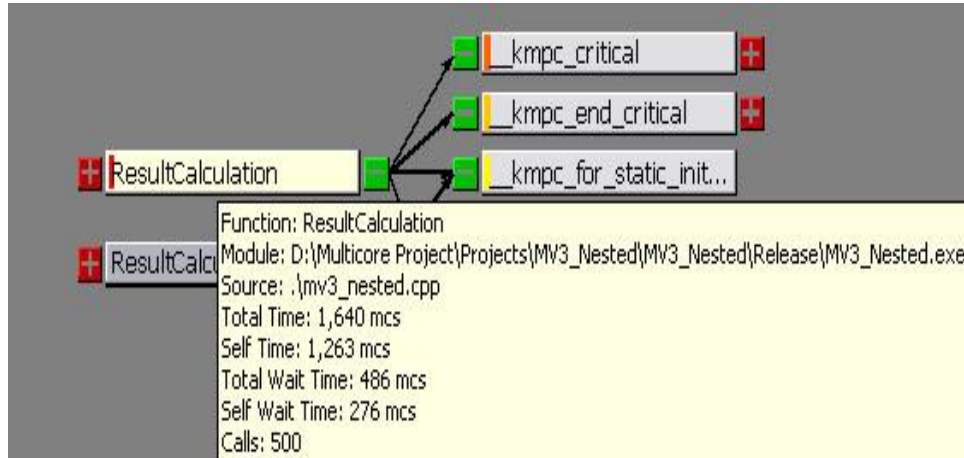


Рис. 6.19. Общее и собственное время выполнения функции *ParallelResultCalculation* при выполнении параллельного алгоритма умножения матрицы на вектор, основанного на блочном разделении данных, реализованного с помощью вложенного параллелизма

Если же использовать информацию о частоте кэш-промахов, то получим:

$$T_p = \left(\frac{n \cdot 2n-1}{p} + \frac{n}{q} \cdot \log_2 q + \frac{n}{q} \right) \cdot \tau + \gamma \cdot n^2 \cdot \left(\alpha + \frac{64}{\beta} \right) + \frac{n}{q} \cdot \delta \quad (6.15)$$

6.7.6. Результаты вычислительных экспериментов

Вычислительные эксперименты для оценки эффективности параллельного алгоритма проводились при тех же условиях, что и ранее выполненные расчеты (см. п. 6.5.5).

1. Первый вариант реализации. Результаты экспериментов приведены в табл. 6.9 и на рис. 6.20.

В табл. 6.10 и на рис. 6.21 представлены результаты сравнения времени выполнения T_p параллельного алгоритма умножения матрицы на вектор с использованием четырех потоков со временем T_p^* , полученным при помощи модели (6.13). Частота кэш-промахов, измеренная с помощью сис-

темы VPS, для четырех потоков значение этой величины была оценена как 0,015.

Таблица 6.9.
Результаты вычислительных экспериментов по исследованию
параллельного алгоритма умножения матрицы на вектор
при блочном разделении данных (четыре потока)

Размер матрицы	Последовательный алгоритм	Параллельный алгоритм	
		Время	Ускорение
1000	0,0076	0,0028	2,7066
2000	0,0303	0,0103	2,9403
3000	0,0688	0,0227	3,0276
4000	0,1222	0,0398	3,0697
5000	0,1909	0,0629	3,0347
6000	0,2748	0,0906	3,0338
7000	0,3741	0,1227	3,0496
8000	0,4894	0,1594	3,0697
9000	0,6186	0,2078	2,9764
10000	0,7637	0,2485	3,0730



Рис. 6.20. Зависимость ускорения от размера матриц при выполнении параллельного алгоритма умножения матрицы на вектор (блочное разбиение матрицы, 4 потока)

Таблица 6.10.

Сравнение экспериментального и теоретического времени выполнения параллельного алгоритма умножения матрицы на вектор, основанного на блочном разбиении матрицы с использованием четырех потоков

Размер матрицы	T_p	T_p^* (calc) (модель)	Модель 6.12 – оценка сверху		Модель 6.13 – уточненная оценка	
			T_p^* (mem)	T_p^*	T_p^* (mem)	T_p^*
1000	0,0028	0,0019	0,0131	0,0150	0,0002	0,0021
2000	0,0103	0,0076	0,0524	0,0600	0,0008	0,0084
3000	0,0227	0,0170	0,1179	0,1349	0,0018	0,0188
4000	0,0398	0,0303	0,2096	0,2399	0,0031	0,0334
5000	0,0629	0,0473	0,3275	0,3748	0,0049	0,0522
6000	0,0906	0,0681	0,4716	0,5397	0,0071	0,0752
7000	0,1227	0,0927	0,6420	0,7346	0,0096	0,1023
8000	0,1594	0,1210	0,8385	0,9595	0,0126	0,1336
9000	0,2078	0,1531	1,0612	1,2144	0,0159	0,1691
10000	0,2485	0,1891	1,3101	1,4992	0,0197	0,2087

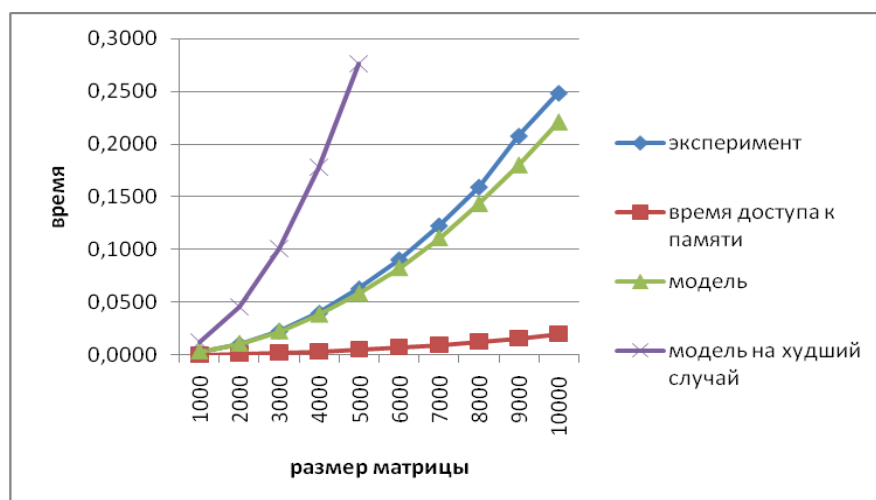


Рис. 6.21. График зависимости экспериментального и теоретического времени выполнения параллельного алгоритма от объема исходных данных при использовании четырех потоков (блочное разбиение матрицы)

2. Второй вариант реализации. Результаты экспериментов приведены в табл. 6.11. и на рис. 6.22.

Таблица 6.11.

Результаты вычислительных экспериментов по исследованию параллельного алгоритма умножения матрицы на вектор при блочном разделении данных и использовании вложенного параллелизма

Размер матрицы	Последовательный алгоритм	Параллельный алгоритм	
		Время	Ускорение
1000	0,0076	0,0050	1,5173
2000	0,0303	0,0147	2,0629
3000	0,0688	0,0296	2,3272
4000	0,1222	0,0488	2,5021
5000	0,1909	0,0720	2,6520
6000	0,2748	0,1016	2,7035
7000	0,3741	0,1343	2,7852
8000	0,4894	0,1739	2,8150
9000	0,6186	0,2144	2,8857
10000	0,7637	0,2632	2,9013

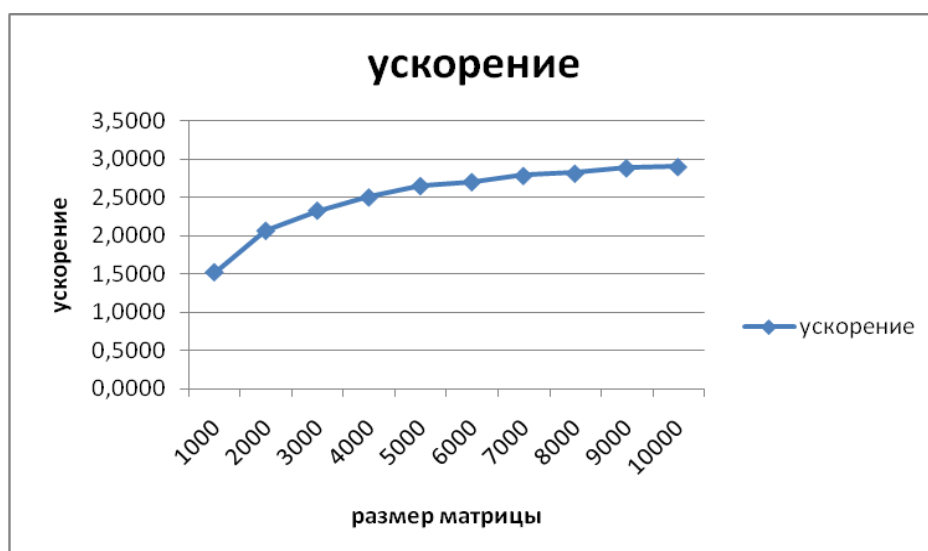


Рис. 6.22. Зависимость ускорения от размера матриц при выполнении параллельного алгоритма умножения матрицы на вектор (блочное разбиение матрицы, 4 потока) с использованием вложенного параллелизма

Как видно из представленных таблиц и графиков, при небольших размерах матрицы и незначительных объемах вычислений алгоритм, основан-

ный на автоматическом распределении вычислительной нагрузки, существенно проигрывает по производительности алгоритму, основанному на «ручном» разделении данных. Это объясняется тем, что в этом случае больший вес имеют накладные расходы на организацию параллелизма. Однако с ростом вычислительной нагрузки доля времени, затраченного на организацию параллелизма, становится меньше, производительность алгоритма повышается. При максимальных размерах матрицы и вектора он практически не уступает алгоритму, реализованному на основе «ручного» разделения.

В табл. 6.12 и на рис. 6.23 представлены результаты сравнения времени выполнения T_p параллельного алгоритма умножения матрицы на вектор с использованием четырех потоков со временем T_p^* , полученным при помощи модели (6.15). Частота кэш-промахов, измеренная с помощью системы VPS, для четырех потоков значение этой величины была оценена как 0,067.

Таблица 6.12.
Сравнение экспериментального и теоретического времени выполнения параллельного алгоритма умножения матрицы на вектор, основанного на блочном разбиении матрицы, с использованием четырех потоков (вложенный параллелизм)

Размер матрицы	T_p	$T_p^* (calc)$ (модель)	Модель 6.14 – оценка сверху		Модель 6.15 – уточненная оценка	
			$T_p^* (mem)$	T_p^*	$T_p^* (mem)$	T_p^*
1000	0,0050	0,0020	0,0131	0,0151	0,0009	0,0029
2000	0,0147	0,0078	0,0524	0,0602	0,0035	0,0113
3000	0,0296	0,0174	0,1179	0,1353	0,0079	0,0253
4000	0,0488	0,0308	0,2096	0,2404	0,0140	0,0448
5000	0,0720	0,0479	0,3275	0,3754	0,0219	0,0698
6000	0,1016	0,0688	0,4716	0,5405	0,0316	0,1004
7000	0,1343	0,0935	0,6420	0,7355	0,0430	0,1365
8000	0,1739	0,1220	0,8385	0,9605	0,0562	0,1782
9000	0,2144	0,1542	1,0612	1,2155	0,0711	0,2253
10000	0,2632	0,1903	1,3101	1,5004	0,0878	0,2781

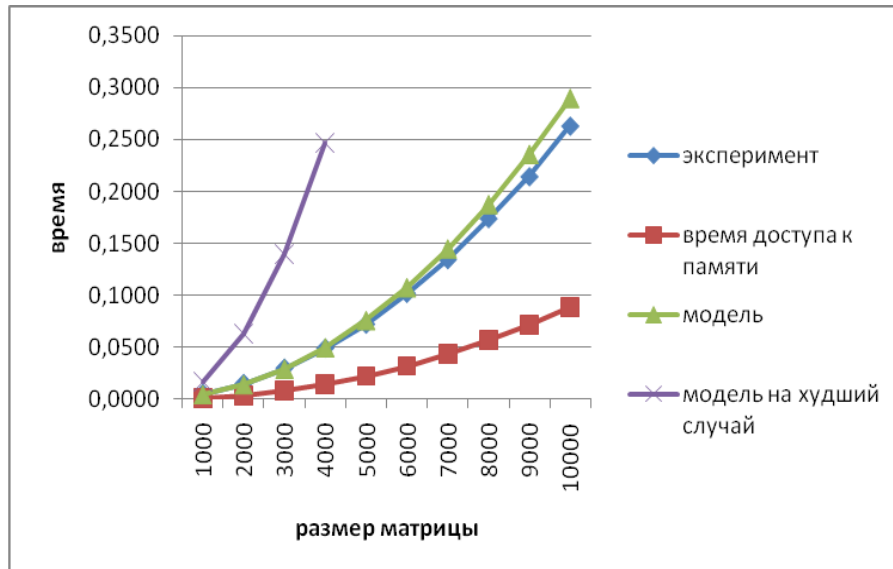


Рис. 6.23. График зависимости экспериментального и теоретического времени выполнения параллельного алгоритма от объема исходных данных при использовании четырех потоков (блочное разбиение, вложенный параллелизм)

6.8. Краткий обзор главы

В данной главе на примере задачи умножения матрицы на вектор рассмотрены возможные схемы разделения матриц между потоками параллельной программы, предназначенной для выполнения на многопроцессорной вычислительной системе с общей памятью и/или на вычислительной системе с многоядерными процессорами. Эти схемы разделения данных являются общими и могут быть использованы для организации параллельных вычислений при выполнении любых матричных операций. В числе излагаемых схем – способы разбиения матриц на *полосы* (по вертикали или горизонтали) или на прямоугольные наборы элементов (*блоки*).

Далее с использованием рассмотренных способов разделения матриц подробно изложены три возможных варианта параллельного выполнения операции умножения матрицы на вектор. Первый алгоритм основан на разделении матрицы между потоками по строкам, второй – на разделении матрицы по столбцам, а третий – на блочном разделении данных. Каждый алгоритм представлен в соответствии с общей схемой, описанной в гл. 3, – вначале определяются базовые подзадачи, затем выделяются информационные зависимости подзадач, далее обсуждается масштабирование и распределение подзадач между вычислительными элементами. В завершение для каждого алгоритма проведен анализ эффективности получаемых па-

раллельных вычислений и приведены результаты вычислительных экспериментов. Для всех рассматриваемых параллельных алгоритмов умножения матрицы на вектор приводятся возможные варианты программной реализации.

Полученные показатели ускорения и эффективности показывают, что все используемые способы разделения данных приводят к равномерной балансировке вычислительной нагрузки, и отличия возникают только в трудоемкости выполняемых информационных взаимодействий между потоками. В этом отношении интересным представляется проследить, как выбор способа разделения данных влияет на характер организации параллельных участков программы, выделить основные различия в использовании потоками общих ресурсов и необходимых операций по синхронизации доступа к ним.

На рис. 6.22 на общем графике представлены показатели ускорения, полученные в результате выполнения вычислительных экспериментов для всех рассмотренных алгоритмов. Как можно заметить, некоторое преимущество по ускорению имеет параллельный алгоритм умножения матрицы на вектор при ленточном разделении по строкам.

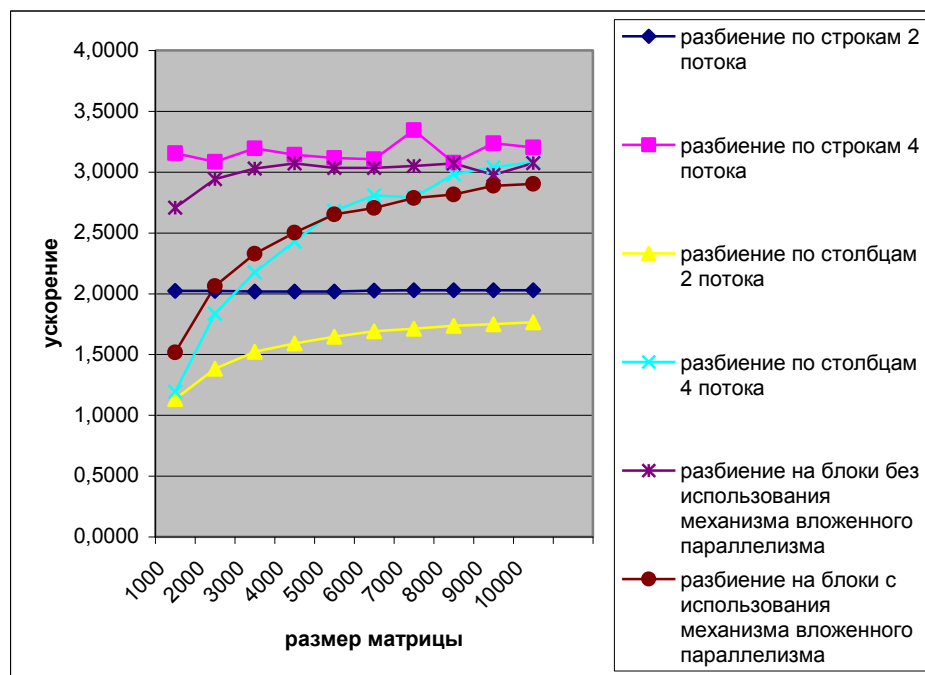


Рис. 6.24. Показатели ускорения рассмотренных параллельных алгоритмов умножения по результатам вычислительных экспериментов

6.9. Обзор литературы

Задача умножения матрицы на вектор часто используется как демонстрационный пример параллельного программирования и, как результат, широко рассматривается в литературе. В качестве дополнительного учебного материала могут быть рекомендованы работы [8,10,72,85]. Широкое обсуждение вопросов параллельного выполнения матричных вычислений выполнено в работе [50].

6.10. Контрольные вопросы

1. Назовите основные способы распределения элементов матрицы между потоками.
2. В чем состоит постановка задачи умножения матрицы на вектор?
3. Какова вычислительная сложность последовательного алгоритма умножения матрицы на вектор?
4. Какие подходы могут быть предложены для разработки параллельных алгоритмов умножения матрицы на вектор?
5. Представьте общие схемы рассмотренных параллельных алгоритмов умножения матрицы на вектор.
6. Проведите анализ и получите показатели эффективности для одного из рассмотренных алгоритмов.
7. Какой из представленных алгоритмов умножения матрицы на вектор обладает лучшими показателями ускорения и эффективности?
8. Может ли использование циклической схемы разделения данных повлиять на время работы каждого из представленных алгоритмов?
9. Какие информационные взаимодействия выполняются для алгоритмов при ленточной схеме разделения данных? В чем различие необходимых операций по организации параллельных участков программы и синхронизации доступа к общим ресурсам при разделении матрицы по строкам и столбцам?
10. Какие информационные взаимодействия выполняются для блочного алгоритма умножения матрицы на вектор?
11. Какие средства технологии OpenMP и функции соответствующей библиотеки оказались необходимыми при программной реализации алгоритмов?

6.11. Задачи и упражнения

1. Выполните реализацию параллельного алгоритма, основанного на ленточном разбиении матрицы на вертикальные полосы. Постройте теоре-

тические оценки времени работы этого алгоритма с учетом параметров используемой вычислительной системы. Проведите вычислительные эксперименты. Сравните результаты реальных экспериментов с ранее подготовленными теоретическими оценками.

2. Выполните реализацию параллельного алгоритма, основанного на разбиении матрицы на блоки. Постройте теоретические оценки времени работы этого алгоритма с учетом параметров используемой вычислительной системы. Проведите вычислительные эксперименты. Сравните результаты реальных экспериментов с ранее подготовленными теоретическими оценками.