

УДК 519.683.4

ОБ ОПТИМИЗАЦИИ ВЫЧИСЛИТЕЛЬНЫХ ПРИЛОЖЕНИЙ ДЛЯ МНОГОПРОЦЕССОРНЫХ СИСТЕМ С ОБЩЕЙ НЕОДНОРОДНОЙ ПАМЯТЬЮ

К. Ю. Богачев¹, А. Р. Миргасимов¹

Для многопроцессорных систем с общей памятью с разным временем доступа процессоров к различным участкам памяти (так называемые NUMA-системы) рассмотрен подход к оптимизации многопоточных приложений, позволяющий максимально использовать вычислительные ресурсы системы при минимальных изменениях в коде приложения. Подход может быть использован в гибридных MPI-многопоточных программах на современных кластерных системах. Приведены результаты численных экспериментов на большом количестве реальных задач. Работа выполнена при поддержке РФФИ (проект № 09-01-00625-а)

Ключевые слова: высокопроизводительные вычисления, гибридные MPI-многопоточные программы, NUMA-системы.

1. Введение. После разработки в 2003 г. процессоров AMD Opteron стали широко доступны системы с разным временем доступа процессоров к различным участкам памяти (так называемые NUMA-системы, Non-Uniform Memory Access). Это привело к значительным изменениям как в операционных системах (планирование заданий, управление памятью) [1], так и в приложениях [2]. С появлением архитектуры процессоров Intel Woodcrest в 2006 г. интерес к NUMA-системам понизился. После создания в 2008 г. процессоров Intel Nehalem стало ясно, что все новые многопроцессорные конфигурации будут NUMA-системами.

Для реализации возможностей NUMA-систем в основном рассматривались две возможности: использование MPI (Message Passing Interface), при котором полностью решается проблема разного времени доступа к памяти ценой пересылки данных из одного процесса, работающего на узле с общей памятью, другому процессу, работающему там же, и использование потоков с разделением данных между ними, при котором требуется значительно менять код традиционного многопоточного приложения.

Ниже рассмотрена третья возможность, основанная на использовании особенностей управления памятью в операционных системах Linux и Windows 7 и не требующая значительного изменения программного кода многопоточного приложения. Проводится сравнение времени работы вычислительного приложения на NUMA-системах на примере решения задачи фильтрации вязкой сжимаемой многофазной смеси в анизотропной пористой среде в реальных нефтегазовых месторождениях.

2. Подход к оптимизации вычислительных приложений для NUMA-систем. Рассмотрим параллельную многопоточную реализацию алгоритма перемножения матрицы на вектор [3]. Скорость работы такой программы во многом зависит от скорости доступа к оперативной памяти. Это делает ее хорошим примером для демонстрации ключевых идей, лежащих в основе метода оптимизации программного кода для NUMA-систем.

В традиционной программной реализации память под матрицы и вектор выделяется в главном потоке. На NUMA-системе это приводит к тому, что физически матрица и вектор располагаются в локальной памяти одного NUMA-узла, что увеличивает как минимум в два раза время доступа к ним потоков, работающих на другом узле.

Покажем, как можно оптимизировать программный код для NUMA-систем, не изменив сам численный алгоритм, и получить ускорение 1.5–2 раза.

Допустим, что порядок матрицы равен $K = MT$, где T — число потоков; пусть матрица хранится в памяти по строкам. Поток с номером i вычисляет результат произведения строк матрицы с номерами $M(i-1)+1, \dots, Mi$ на вектор. Будем говорить, что эти строки матрицы принадлежат i -му потоку. Таким образом, матрица делится на T непрерывных частей, каждая из которых обрабатывается собственным потоком.

¹ Московский государственный университет им. М. В. Ломоносова, механико-математический факультет, Ленинские горы, 119899, Москва; К. Ю. Богачев, доцент, e-mail: bogachev@mech.math.msu.su; А. Р. Миргасимов, аспирант, e-mail: mirgasimov.almaz@gmail.com

Ключевая идея оптимизации заключается в реализации двух условий, которые позволят минимизировать число обращений к нелокальным участкам памяти:

- каждый из вычислительных потоков работает на фиксированном узле NUMA-системы,
- те части матрицы и вектора, с которыми вычислительный поток оперирует большую часть времени, физически находятся в локальной памяти этого узла.

Для выполнения первого условия в каждом из вычислительных потоков сразу после его создания сделаем системный вызов, который привяжет этот i -й поток к i -му логическому процессору.

Для удовлетворения второго условия с минимальным переписыванием программы воспользуемся особенностью механизма выделения памяти в операционных системах, поддерживающих NUMA. В системах Linux и Windows 7 при выделении памяти посредством функции malloc отображение затребованной виртуальной памяти на физическую строится не сразу. Когда поток обращается к блоку виртуальной памяти, для которого еще не построено отображение, происходит страничный промах (page fault), который инициирует создание отображения запрошенного блока виртуальной памяти на локальную память того узла NUMA-системы, с которого был произведен запрос.

При традиционном подходе выделения памяти в многопоточных приложениях главный поток запрашивает память у операционной системы и инициализирует ее, скажем, нулями. На NUMA-системе это приводит к тому, что вся запрошенная память выделится в локальной памяти узла, на котором работал главный поток. Остальные потоки в это время простаивают, дожидаясь завершения этой операции.

Модифицируем указанный подход, чтобы его можно было эффективно использовать на NUMA-системе. Будем выполнять операцию выделения памяти в два этапа. На первом этапе главный поток запрашивает у операционной системы память требуемого объема. Затем следует точка синхронизации, после которой указатель на затребованную память можно использовать во всех потоках. На втором этапе выделенная память параллельно инициализируется всеми потоками, причем каждый из потоков инициализирует только свою часть запрошенной памяти.

Предложенный метод инициализации памяти в два этапа эффективен по двум причинам. Во-первых, та часть памяти, к которой поток будет обращаться чаще всего, выделяется в локальной памяти узла NUMA-системы, на котором тот работает. Во-вторых, ускоряется операция ее инициализации за счет параллелизма в работе памяти. Минусом этого метода является наличие дополнительной точки синхронизации потоков.

Эти модификации затрагивают ту часть программного кода, которую можно назвать подготовительной стадией алгоритма — создание вычислительных потоков, заполнение матрицы и вектора. Тем самым алгоритмически существенная часть программного кода, ответственная за непосредственное перемножение матрицы на вектор, остается неизменной.

На основе этих идей представленный метод оптимизации кода для NUMA-систем можно применить к любой многопоточной программе, в которой вычислительные потоки большую часть времени проводят, обращаясь к своей части всей выделенной процессу памяти. Получаемый выигрыш зависит от доли операций с памятью в общем времени работы программы.

3. Постановка задачи фильтрации. Рассмотрим следующие стандартные уравнения многофазной изотермической модели черной нефти в трехмерной области, использованные также в работах [5, 6]:

$$\begin{aligned} \frac{\partial}{\partial t} (\phi N_c) &= \operatorname{div} \sum_{l=o,w,g} x_{cl} \xi_l \left(\mathbf{k} \frac{k_{rl}}{\mu_l} (\nabla p_l - \gamma_l \nabla \mathbf{D}) \right) + q_c, \quad c = 1, \dots, n_c, \\ p_o - p_g &= P_{cog}, \quad p_o - p_w = P_{c ow}, \\ S_w + S_o + S_g &= 1, \\ N_w &= \frac{S_w}{B_w}, \quad N_o = \frac{S_o}{B_o}, \quad N_g = \frac{S_g}{B_g} + R_{g,o} \frac{S_o}{B_o}. \end{aligned} \tag{1}$$

Здесь n_c — количество компонентов в смеси; $R_{g,o} = R_{g,o}(p_o)$ — растворимость газа в нефтяной фазе; $B_l = B_l(p_l)$ — коэффициент объемного расширения фазы l , $l = o$ (oil), g (gas), w (water); пористость среды $\phi = \phi(p_w, p_o, p_g, x, y, z)$; $x_{cl} = x_{cl}(p_w, p_o, p_g, \mathbf{N})$ — молярная доля компонента c , $c = 1, \dots, n_c$ в фазе l ; $\mathbf{N} = (N_1, \dots, N_{n_c})$ — вектор молярных плотностей; $\xi_l = \xi_l(p_l, \mathbf{N})$ — молярная плотность фазы l ; $\mathbf{k} = \mathbf{k}(p_w, p_o, p_g, x, y, z)$ — тензор абсолютной проницаемости; $k_{rl} = k_{rl}(S_w, S_g)$ — относительная проницаемость фазы l ; $\mu_l = \mu_l(p_l)$ — вязкость фазы l ; $\gamma_l = \rho_l g$ — вертикальный градиент гидростатического давления в фазе l ; \mathbf{g} — гравитационная постоянная; $\mathbf{D} = \mathbf{D}(x, y, z)$ — вектор глубины (сверху вниз); $\rho_l = \rho_l(p_l)$ — массовая плотность фазы l ; $P_{cog} = P_{cog}(S_g)$ — капиллярное давление в системе нефть–газ; $P_{c ow} = P_{c ow}(S_w)$ — капиллярное давление в системе вода–нефть; $q_c = q_c(p_w, p_o, p_g, \mathbf{N}, t, x, y, z)$ — источник

компонента c (скважина).

Неизвестными в этой системе являются

1) $N_c = N_c(t, x, y, z)$ — молярная плотность компонента c (для модели черной нефти компонентами служат вода, нефть и газ);

2) $S_l = S_l(t, x, y, z)$ — насыщенность фазы l ;

3) $p_l = p_l(t, x, y, z)$ — давление в фазе l .

В самом распространенном случае количество компонентов n_c равно трем. Для этой системы уравнений задаются начальные условия, а также на внешней границе резервуара ставятся условия непротекания.

4. Практическое применение подхода к оптимизации вычислительных приложений для NUMA-систем. Продемонстрируем практическое применение оптимизации кода для NUMA-систем на примере решения задачи (1) фильтрации вязкой сжимаемой многофазной смеси в анизотропной пористой среде, реализованного в гидродинамическом пакете моделирования tNavigator.

Для аппроксимации уравнений по времени используется полностью неявная схема. Аппроксимацией по пространственным переменным методом конечных объемов на неравномерной сетке первоначальная задача сначала сводится к системе нелинейных алгебраических уравнений вида

$$F(\mathbf{p}, N_1, \dots, N_{n_c}) = 0,$$

где $\mathbf{p} = (p^i)$, $\mathbf{N}_c = (N_c^i)$ — векторы значений давления и молярных плотностей в блоках сетки. Для решения системы нелинейных уравнений $F(\mathbf{U}) = 0$, $\mathbf{U} \equiv (\mathbf{p}, \mathbf{N})$ используется метод Ньютона

$$\mathbf{U}^{m+1} = \mathbf{U}^m - \left(\frac{\partial F(\mathbf{U}^m)}{\partial \mathbf{U}} \right)^{-1} F(\mathbf{U}^m).$$

Здесь $\frac{\partial F(\mathbf{U}^m)}{\partial \mathbf{U}}$ — матрица $R^{n_c * K} \times R^{n_c * K}$, где K — число блоков сетки. На каждом шаге метода Ньютона решается линейная система с несимметричной матрицей $\frac{\partial F(\mathbf{U}^m)}{\partial \mathbf{U}}$, т.е. задача сводится к решению системы линейных алгебраических уравнений

$$A\mathbf{x} = \mathbf{r}, \tag{2}$$

в которой матрица A — якобиан из метода Ньютона.

Построение системы линейных уравнений и ее решение занимают большую часть времени работы пакета tNavigator. Именно эти два этапа были подвергнуты оптимизации для NUMA-систем.

Перед тем как приступить к оптимизации параллельного многопоточного кода для NUMA-систем, необходимо определить, какие объекты, имеющие размер порядка K (назовем их большими), используются в программе и можно ли каждый из этих объектов разбить на части, которые в основном обрабатываются отдельными потоками.

Первое множество больших объектов (обозначим его через T_1) включает в себя различные свойства фаз, компонентов и породы в каждом блоке сетки. Во-первых, это главные переменные решаемой задачи — молярные плотности компонентов и давление в фазе “нефть”. Во-вторых, все, что так или иначе входит в коэффициенты уравнения (1) — относительные проницаемости, вязкости, насыщенности фаз, пористость, проводимость, коэффициенты масштабирования фазовых проницаемостей и т.д. Этот список на практике получается гораздо длиннее приведенного и зависит от множества физических явлений, учитываемых в математической модели. В-третьих, это производные перечисленных свойств по главным переменным.

Указанные данные, в основном, используются для вычисления коэффициентов системы линейных уравнений (2). Это процедуру можно представить следующим циклом: для каждого блока $i = 1, \dots, K$, используя данные i -го блока и геометрически соседних с ним, посчитать коэффициенты линейной системы, описывающих физику процесса в i -м блоке.

В многопоточном приложении множество всех блоков разбивается между M потоками — i -й поток обрабатывает блоки с номерами $K_{i-1}, \dots, K_i - 1$, где $1 = K_0 < K_1 < \dots < K_M = K$. Как правило, разбиение получается таким, что геометрически соседние блоки обрабатываются одним потоком. Таким образом, i -й поток при вычислении своей части коэффициентов линейной системы, в основном, будет обращаться к элементам в позициях $K_{i-1}, \dots, K_i - 1$ массивов размера K из множества T_1 .

Система линейных уравнений (2) решается методом BCGS (Biconjugate Gradient Stabilized) [7] с использованием предобусловливателя класса ILU (Incomplete LU-decomposition) [4]. На этом этапе численные операции выполняются на другом множестве больших объектов, обозначим его через T_2 . В него входят

Таблица 1

Тестовая система 1: процессор Intel Xeon E5520 2.27 GHz, число ядер = 4, число процессоров = 2, оперативная память DDR-III 1067 MHz 12 Gb		
а	Умножение матрицы $K \times K$ на вектор длины K	Размерность задачи $K = 20000$ Объем использованной памяти — 3052 Мб
б	Решение системы линейных уравнений методом BCGS	Размер матрицы $K = 1094421$ Число ненулевых элементов в матрице = 7484949 Объем использованной памяти — 702 Мб
в	Расчет реального месторождения	Размерность задачи (число блоков в модели) $K = 1094421$ Объем использованной памяти — 1921 Мб Время расчета с использованием одного потока — 1 час 2 минуты

Таблица 2

Тестовая система 2: процессор Intel Xeon X7560 2.27GHz, число ядер = 8, число процессоров = 4, оперативная память DDR-III 1067MHz 128Gb	
Расчет реального месторождения	Размерность задачи (число блоков в модели) $K = 2418989$ Объем использованной памяти — 6133 Мб Время расчета с использованием одного потока — 26 часов 47 минут

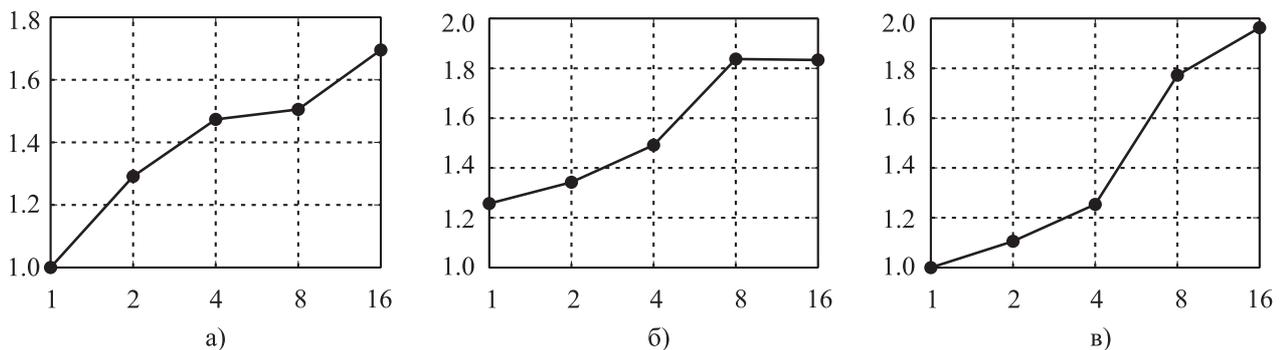


Рис. 1. Графики ускорения работы программ после применения оптимизаций для NUMA-систем в зависимости от числа потоков: а) умножение матрицы на вектор, б) решение системы линейных уравнений, в) расчет реального месторождения

матрица A , хранящаяся в MSR-формате, матрица предобусловливателя для каждого потока, располагающаяся в выделенной ему памяти, и некоторое количество векторов длины K , используемых в методе BCGS. При решении системы i -й поток, в основном, оперирует со строками $K_{i-1}, \dots, K_i - 1$ матрицы A , с элементами в тех же позициях векторов длины K и со своей матрицей предобусловливателя.

Отмеченное свойство вычислительных потоков обращаться, в основном, только к своим частям массивов дает возможность разбивать большие объекты на части. Таким образом, теперь оптимизация программного кода для NUMA-систем становится такой же, как и в случае перемножения матрицы на вектор:

- использование системного вызова, привязывающего i -й поток к i -му логическому процессору;
- использование двухэтапной инициализации памяти под массивы из множеств T_1 и T_2 , в которой каждый из потоков инициализирует только свою часть этих массивов.

5. Численные эксперименты. Описанный подход был реализован в составе гидродинамического пакета tNavigator и протестирован на параллельной ЭВМ с двумя процессорами Intel Nehalem и на параллельной ЭВМ с четырьмя процессорами Intel Nehalem EX. Для тестирования была выбрана гидродинамическая модель реального месторождения, содержащая несколько миллионов блоков. Тестовая система 1 представлена в табл. 1, тестовая система 2 — в табл. 2.

Графики ускорения работы программ после применения оптимизаций для NUMA-систем показаны на рис. 1а–1в.

На рис. 2а показан график ускорения работы программ после применения оптимизаций для NUMA-систем, на рис. 2б — график масштабируемости параллельного расчета.

Оптимизация программы для NUMA-системы позволила достичь ускорения в 21 раз на системе с 32 физическими ядрами, в то время как ускорение неоптимизированной версии равно 11.

6. Заключение. Предложен подход, позволяющий значительно сократить время на разработку эффективных вычислительных приложений для NUMA-систем. Результаты численных экспериментов показывают, что при использовании всех физических ядер параллельной ЭВМ оптимизированное для NUMA-систем приложение работает в 1.5 – 2 раза быстрее, чем его исходная версия.

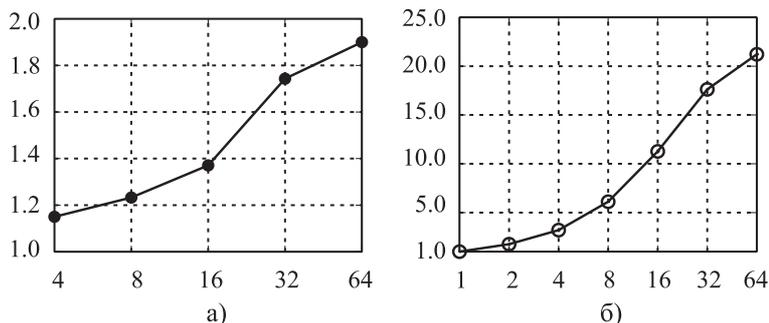


Рис. 2. Расчет реального месторождения, на горизонтальной оси отложено число потоков

СПИСОК ЛИТЕРАТУРЫ

1. Lameter C. Local and remote memory: memory in a Linux/NUMA system (ftp://ftp.kernel.org/pub/linux/kernel/people/christoph/pmig/numamemory.pdf, 2006).
2. David E.Ott. Optimizing software applications for NUMA. DDJ Magazine, 2009.
3. Богачев К.Ю. Основы параллельного программирования. М.: Бинум, 2003.
4. Богачев К.Ю., Жаблицкий Я.В. Блочные предобусловливатели класса ILU для задач фильтрации многокомпонентной смеси в пористой среде // Вестн. Моск. ун-та. Матем. Механ. 2009. № 5. 19–25.
5. Богачев К.Ю., Горелов И.Г. Применение параллельного предобусловливателя CPR к задаче фильтрации вязкой сжимаемой жидкости в пористой среде // Вычислительные методы и программирование. 2008. 9, № 2. 35–41.
6. Богачев К.Ю., Мельниченко Н.С. О пространственной аппроксимации методом подсеток для задачи фильтрации вязкой сжимаемой жидкости в пористой среде // Вычислительные методы и программирование. 2008. 9, № 2. 42–50.
7. Saad Y. Iterative methods for sparse linear systems. Philadelphia: SIAM, 2003.

Поступила в редакцию
21.05.2010