

УДК 004.272

М.Г. Курносов

ОПТИМИЗАЦИЯ АЛГОРИТМОВ КОЛЛЕКТИВНЫХ ОБМЕНОВ ИНФОРМАЦИЕЙ МЕЖДУ ВЕТВЯМИ ПАРАЛЛЕЛЬНЫХ ПРОГРАММ В РАСПРЕДЕЛЕННЫХ ВЫЧИСЛИТЕЛЬНЫХ СИСТЕМАХ С ИЕРАРХИЧЕСКОЙ СТРУКТУРОЙ¹

Предлагается метод оптимизации алгоритмов реализации коллективных обменов информацией между ветвями параллельных программ в иерархических распределенных вычислительных системах (ВС). Метод поясняется на примере создания алгоритмов трансляционно-циклических обменов (ТЦО, All-to-all Broadcast), учитывающих иерархическую структуру распределенных ВС. Показана эффективность разработанных алгоритмов ТЦО в сравнении с существующими аналогами. Накладные расходы на оптимизацию незначительны и компенсируются сокращением времени реализации ТЦО предложенными алгоритмами.

***Ключевые слова:** операции коллективных обменов информацией, модель передачи сообщений, параллельное программирование, распределенные вычислительные системы.*

Начиная с середины 1960-х годов распределенные вычислительные системы (ВС) активно используются как инструментальные средства решения сложных задач в различных отраслях науки и техники. В архитектурном плане распределенная ВС [1, 2] представляется композицией множества взаимодействующих элементарных машин (оснащенных средствами коммуникаций и внешними устройствами) и сети межмашинных связей. Элементарная машина (ЭМ) – это основной функциональный и структурный элемент ВС; конфигурация ЭМ допускает варьирование в широких пределах – от процессорного ядра до ЭВМ. Все основные ресурсы распределенных ВС (арифметико-логические устройства, память, средства управления и коммуникаций) являются логически и технически рассредоточенными. Число ЭМ в распределенных ВС допускает варьирование от нескольких единиц до сотен тысяч.

Параллельные алгоритмы и программы для распределенных ВС преимущественно разрабатываются в модели передачи сообщений (Message Passing). В этой модели ветви параллельной программы взаимодействуют друг с другом путем обменов информационными сообщениями по каналам межмашинных связей ВС. Выделяют два типа обменов информацией между ветвями параллельных программ [1, 2]: дифференцированный (point-to-point) и коллективные (collective) обмены. При дифференцированном обмене осуществляется передача информации из одной ветви в любую другую ветвь. Коллективные обмены подразделяются на несколько видов: трансляционный (ТО, One-to-all Broadcast), трансляционно-циклический (ТЦО, All-to-all Broadcast) и коллекторный обмены (КО, All-to-one

¹ Работа выполнена в рамках интеграционного проекта № 113 СО РАН, при поддержке РФФИ (гранты № 08-07-00018, 09-07-90403, 09-07-13534, 09-02-12106), Совета по грантам Президента РФ для поддержки ведущих научных школ (грант НШ-5176.2010.9) и в рамках государственного контракта № 02.740.11.0006 с Минобрнауки РФ.

Reduce). При трансляционном обмене данные из одной ветви передаются во все остальные; при трансляционно-циклическом обмене информация из ветвей передается каждой ветви и принимается из всех; коллекторный обмен подразумевает прием информации из всех ветвей в одну.

Анализ использования в параллельных алгоритмах и программах схем обменов информацией показывает, что до 80 % времени обменов приходится на коллективные операции [3]. Как правило, количество обращений в алгоритмах и программах к коллективным операциям обменов имеет функциональную зависимость от размера входных данных и в среднем находится в интервале от 10^1 до 10^4 .

В настоящее время в коммуникационных библиотеках стандарта MPI и системах параллельного программирования (в частности, в модели PGAS – Partitioned Global Address Space) для реализации коллективных обменов используются алгоритмы рассылки данных по кольцу, рекурсивного сдваивания, алгоритм Дж. Брука (J. Bruck) и алгоритмы, упорядочивающие ветви в деревья различных видов [4]. Перечисленные алгоритмы характеризуются различным временем выполнения и опираются на предположение об однородности каналов связи между ЭМ распределенных ВС. Однако современные системы являются мультиархитектурными, для них характерны иерархическая структура и зависимость времени передачи данных между ЭМ от их размещения в системе [5 – 7].

В данной работе предлагается метод оптимизации алгоритмов реализации коллективных обменов, учитывающий иерархическую структуру современных распределенных ВС и производительность каналов связи между ЭМ системы. Суть подхода демонстрируется на примере создания структурно-ориентированных алгоритмов реализации трансляционно-циклических обменов.

1. Алгоритмы реализации трансляционно-циклических обменов

При реализации ТЦО каждая ветвь $i \in \{0, 1, \dots, n-1\}$ параллельной программы располагает локальным сообщением a_i размером m байт. Требуется отправить сообщение a_i всем ветвям и получить сообщения от каждой. По окончании ТЦО все ветви должны содержать в своей памяти n сообщений a_0, a_1, \dots, a_{n-1} , упорядоченных по номерам ветвей, которым они принадлежат.

Кольцевой алгоритм (Ring) реализации ТЦО организует параллельные ветви в логическое кольцо [4]. На шаге $k = 0, 1, \dots, n-2$ ветвь i передает сообщение $a_{(i-k+n) \bmod n}$ ветви $(i+1) \bmod n$ и принимает сообщение $a_{(i-k-1) \bmod n}$ от ветви $(i-1+n) \bmod n$.

Алгоритм рекурсивного сдваивания (Recursive Doubling) допускает реализацию ТЦО только для случая n равно степени двойки. На шаге $k = 0, 1, \dots, \log_2 n - 1$ ветви i и $i \oplus 2^k$ обмениваются ранее принятыми 2^k сообщениями (здесь и далее \oplus – сложение по модулю 2). Таким образом, на шаге 0 передается сообщение размером m байт, на шаге 1 – размером $2m$ байт и т. д.

Алгоритм Дж. Брука (J. Bruck) [4] подобен алгоритму рекурсивного сдваивания, но допускает реализацию ТЦО для произвольного количества n параллельных ветвей. На шаге $k = 0, 1, \dots, \lceil \log_2 n \rceil - 1$ ветвь i передает все принятые сообщения ветви $(i - 2^k + n) \bmod n$ и принимает сообщения от ветви $(i + 2^k) \bmod n$. Сообщения размещаются в памяти со смещением, поэтому в конце работы алгоритма каждая ветвь i циклически сдвигает сообщения на i позиций вниз.

Параллельная программа, реализующая ТЦО, может быть представлена информационным графом $G = (V, E)$, вершинам которого соответствуют ветви про-

граммы, а ребрам – информационные обмены между ними. Вес d_{ij} ребра в графе отражает объем данных, переданных по нему при реализации алгоритма. Объемы данных, передаваемых между ветвями параллельной программы при реализации ТЦО алгоритмами рекурсивного сдвигания и Дж. Брука, различны. По этой причине время реализации ТЦО рассматриваемыми алгоритмами в иерархических распределенных ВС зависит от того, как ветви параллельной программы распределены по элементарным машинам системы. Рассмотрим модель иерархической организации распределенных ВС.

2. Модель распределенных ВС с иерархической структурой

Распределенная вычислительная система с иерархической структурой, укомплектованная N однородными ЭМ, может быть представлена в виде дерева, содержащего L уровней [6, 7]. Каждый уровень системы образован отдельным видом функциональных модулей (например, телекоммуникационные шкафы, вычислительные узлы и т. п.), которые объединены каналами связи своего уровня. Введем следующие обозначения: $t_l(m)$ – время передачи сообщения размером m байт между парой элементарных машин ВС через каналы связи уровня $l \in \{1, 2, \dots, L\}$; b_l – максимальное значение пропускной способности каналов связи на уровне уровня l ; $z(p, q)$ – номер уровня функционального модуля, являющегося ближайшим общим предком для ЭМ $p, q \in \{1, 2, \dots, N\}$, иначе говоря, номер уровня, через который они взаимодействуют.

Вычислительные системы, представленные в списке TOP500 (35-я редакция), имеют как минимум два уровня в иерархической организации – сеть связи между вычислительными узлами и их общую память.

На рис. 1 показано время передачи данных через каналы связи различных уровней между ЭМ кластерных ВС: кластера Центра параллельных вычислительных технологий ГОУ ВПО «Сибирский государственный университет телекоммуникаций и информатики» (ЦПВТ ГОУ ВПО «СибГУТИ») и кластера Информационно-вычислительного центра ГОУ ВПО «Новосибирский государственный университет» (ИВЦ ГОУ ВПО «НГУ»).

Время передачи данных между ЭМ в иерархических ВС сокращается с уменьшением номера уровня, через каналы связи которого они взаимодействуют. В то же время, в зависимости от конфигурации распределенной ВС, могут существовать интервалы размеров сообщений, для которых последнее утверждение не выполняется. Например (рис. 1, б), для сообщений с размерами $m > 2^{22}$ байт. Рассмотрим особенности реализации ТЦО в иерархических распределенных ВС.

3. Реализации ТЦО в иерархических распределенных ВС

Обозначим через $x_i \in \{1, 2, \dots, N\}$ номер ЭМ, на которую распределена ветвь $i \in \{0, 1, \dots, n-1\}$ программы; h – количество функциональных модулей уровня L (вычислительных узлов), выделенных для реализации программы; q_1, q_2, \dots, q_h – номера выделенных функциональных модулей; s_r – количество ЭМ функционального модуля q_r , выделенных для реализации программы, $r \in \{1, 2, \dots, h\}$.

На рис. 2 приведен пример реализации ТЦО алгоритмом Дж. Брука на вычислительном кластере с иерархической организацией. Время выполнения алгоритма с распределением ветвей по ЭМ стандартными средствами библиотеки MPI (MPICH2 1.2.1) составило 567 мкс, а время реализации с учетом графа алгоритма и иерархической организации кластерной ВС – 324 мкс.

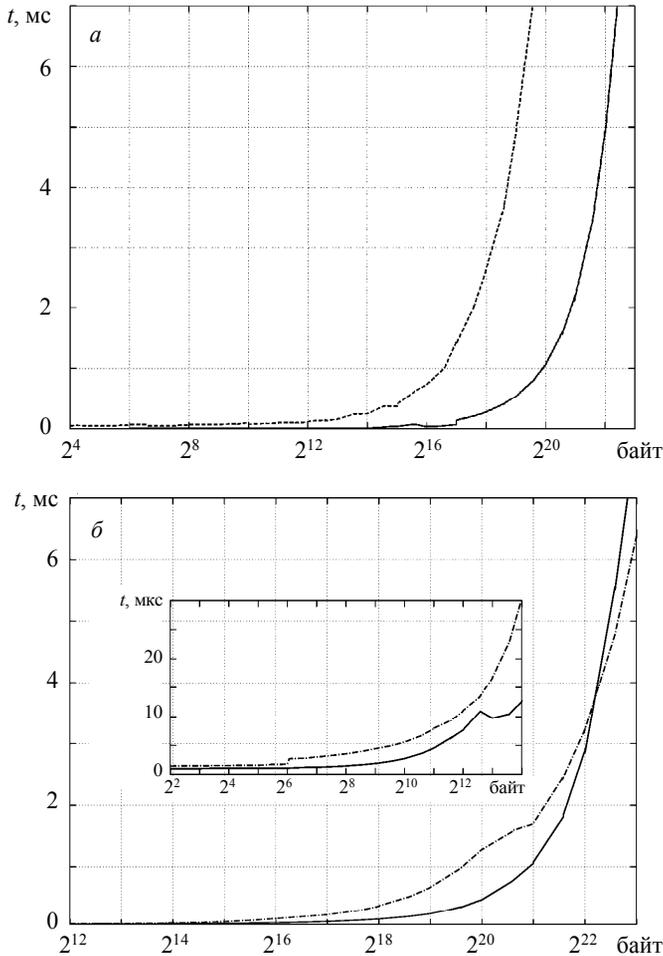


Рис. 1. Время передачи данных между элементарными машинами кластерной ВС через общую память вычислительного узла и сеть межузловых связей: *a* – кластер ЦПВТ ГОУ ВПО «СибГУТИ»; *б* – кластер ИВЦ ГОУ ВПО «НГУ»; — — — — разделяемая память; - · - · - сеть Gigabit Ethernet; - · - · - сеть InfiniBand 4x DDR

Причиной сокращения времени выполнения ТЦО в 1,75 раз (рис. 2, б) является распределение ветвей, обменивающихся большими объемами данных, на один вычислительный узел, в рамках которого они взаимодействуют через его общую память. Учитывая последнее, разработаны метод и алгоритмы реализации ТЦО в иерархических распределенных ВС.

Накладные расходы на реализацию ТЦО с заданным распределением $X = (x_0, x_1, \dots, x_{n-1})$ параллельных ветвей по ЭМ системы и заданным информационным графом $G = (V, E)$ алгоритма могут быть выражены как сумма $T(X)$ времен передачи данных между ветвями программы

$$T(X) = \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} d_{ij} / b_{z(x_i, x_j)},$$

где d_{ij} – вес ребра в графе алгоритма, а $b_{z(x_i, x_j)}$ – пропускная способность каналов связи на уровне $z(x_i, x_j)$, через которые взаимодействуют ветви i и j .

В основе разработанного метода лежит процедура распределения параллельных ветвей программы, обменивающихся большими объемами данных при реализации ТЦО заданным алгоритмом, на ЭМ одного функционального модуля. Это обеспечивает локализацию взаимодействий ветвей через каналы связи функционального модуля и, как следствие, сокращение времени информационных обменов. Например, при реализации ТЦО алгоритмом Дж. Брука суммарный объем данных, передаваемых между ветвями 0, 2, 4 и 6, больше объема данных, отправляемых другим ветвям, поэтому распределение этих ветвей на один вычислительный узел (рис. 2, б) обеспечило их взаимодействие через общую память узла и сокращение времени реализации ТЦО.

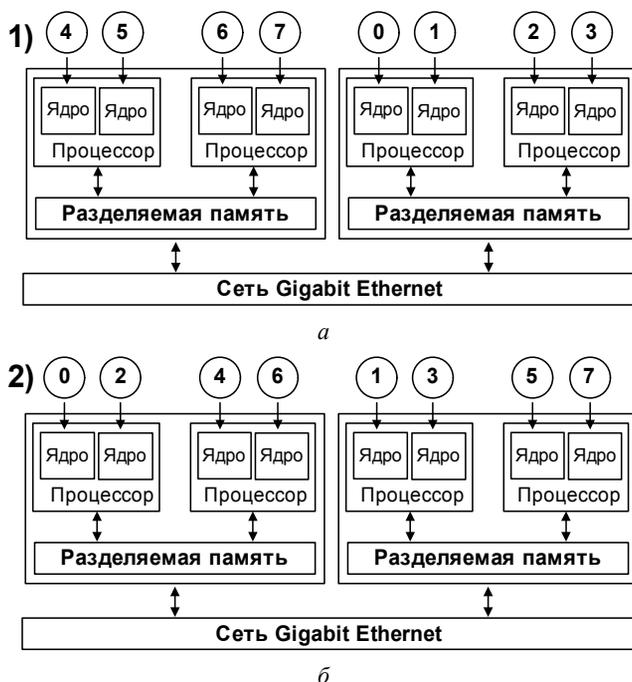


Рис. 2. Реализация ТЦО алгоритмом Дж. Брука на кластерной ВС ($n = 8$; $m = 2048$; $L = 2$): а – реализация ТЦО с распределением ветвей стандартным алгоритмом библиотеки MPI; б – реализация ТЦО с учетом иерархической организации ВС

Метод подразумевает (суб)оптимальное распределение ветвей по ЭМ и реализацию ТЦО в соответствии с ним. Рассмотрим основные шаги метода, осуществляемые до реализации ТЦО.

1. Формируется взвешенный информационный граф $G = (V, E)$ алгоритма реализации ТЦО для $m = 1$ с количеством вершин n .

2. Строится разбиение $R = (r_0, r_1, \dots, r_{n-1})$ информационного графа $G = (V, E)$ на h непересекающихся подмножеств V_1, V_2, \dots, V_h с целью минимизации суммарного веса рёбер, инцидентных различным подмножествам разбиения. Через $r_i \in \{1, 2, \dots, h\}$ обозначен номер подмножества разбиения, к которому отнесена

вершина $i \in V$. Количество элементов в подмножестве V_u должно быть равно заданному числу $s_u, u = 1, 2, \dots, h$. Параметры h и s_u определены ранее.

3. Ветвям с номерами из подмножества V_u назначаются номера ветвей, распределенных на элементарные машины функционального модуля $q_u, u = 1, 2, \dots, h$. Таким образом, строится взаимно однозначное отображение $\pi(i)$, которое сопоставляет исходным номерам $0, 1, \dots, i, \dots, n-1$ ветвей новые номера $\pi(0), \pi(1), \dots, \pi(n-1)$. Через $\pi^{-1}(i)$ обозначим отображение, обратное к $\pi(i)$.

На рис. 3 показано применение описанного метода для алгоритма Дж. Брука на кластерной ВС с иерархической структурой.

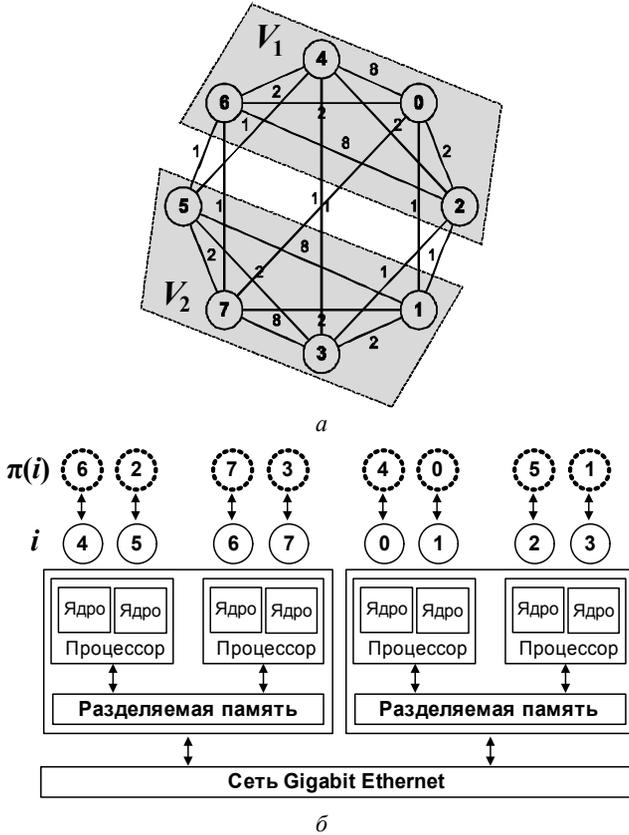


Рис. 3. Пример оптимизации алгоритма Дж. Брука ($L = 2$; $n = 8$; $h = 2$; $s_1 = 4$; $s_2 = 4$): а – приближенное разбиение графа алгоритма Дж. Брука на подмножества V_1 и V_2 ; б – отображение $\pi(i)$ номеров ветвей программы

Алгоритмы рекурсивного сдваивания и Дж. Брука требуют модификации для реализации ТЦО с заданным распределением (отображениями $\pi(i)$ и $\pi^{-1}(i)$) ветвей по ЭМ. В исходных алгоритмах каждая ветвь i располагает сообщением a_i , но в соответствии с полученным распределением ветвь должна осуществлять ТЦО под новым номером $\pi(i)$. Предложим версии алгоритмов рекурсивного сдваивания и Дж. Брука, обеспечивающие реализацию ТЦО с заданными распределениями ветвей по ЭМ.

Алгоритм Bruck exch. Ветвь i передает свое сообщение a_i ветви $\pi^{-1}(i)$, принимает от ветви $\pi(i)$ сообщение и делает его начальным. На шаге $k = 0, 1, \dots, \lceil \log_2 n \rceil - 1$ ветвь i передает все принятые сообщения ветви $\pi^{-1}((i' - 2^k + n) \bmod n)$ и принимает сообщения от ветви $\pi^{-1}((i' + 2^k) \bmod n)$, где $i' = \pi(i)$. Далее каждая ветвь циклически сдвигает сообщения вниз на i' позиций.

Алгоритм Bruck reorder. На шаге $k = 0, 1, \dots, \lceil \log_2 n \rceil - 1$ ветвь i передает все принятые сообщения ветви $\pi^{-1}((i' - 2^k + n) \bmod n)$ и принимает сообщения от ветви $\pi^{-1}((i' + 2^k) \bmod n)$, где $i' = \pi(i)$. Далее каждая ветвь переставляет сообщение из позиции $j = 0, 1, \dots, n - 1$ в позицию $\pi^{-1}((j + i') \bmod n)$.

Алгоритм recursive doubling exch. Ветвь i передает свое сообщение a_i ветви $\pi^{-1}(i)$, принимает от ветви $\pi(i)$ сообщение и делает его начальным. На шаге $k = 0, 1, \dots, \log_2 n - 1$ ветвь i и $\pi^{-1}(i \oplus 2^k)$ обмениваются принятыми 2^k сообщениями.

Алгоритм recursive doubling reorder. Ветвь i сдвигает свое сообщение a_i из позиции i в позицию $\pi(i)$. На шаге $k = 0, 1, \dots, \log_2 n - 1$ ветвь i и $\pi^{-1}(i \oplus 2^k)$ обмениваются ранее принятыми 2^k сообщениями. Далее каждая ветвь переставляет сообщение из позиции $j = 0, 1, \dots, n - 1$ в позицию $\pi^{-1}(j)$.

Формирование, разбиение графа и построение отображений $\pi(i)$ и $\pi^{-1}(i)$ осуществляется единожды при создании подсистемы ЭМ (например, при создании коммутаторов в библиотеках стандарта MPI). После этого построенные отображения многократно используются для реализации ТЦО. Все шаги метода эффективно реализуемы в параллельном виде на подсистеме ЭМ, выделенной для реализации программы. Также допустима организация распределенного хранения информационного графа на ЭМ системы. На протяжении всего времени выполнения программы ветвям требуется лишь информация об отображении смежных с ними ветвей. Последнее, для алгоритмов рекурсивного сдваивания и Дж. Брука, требует хранения каждой ветвью в памяти порядка $2 \lceil \log_2 n \rceil$ байт. Сказанное выше обеспечивает применимость метода для большемасштабных распределенных ВС.

4. Экспериментальное исследование алгоритмов

Разработанный метод реализован в коммуникационной библиотеке ToroMPI, созданной и развиваемой ЦПВТ ГОУ ВПО «СибГУТИ» совместно с лабораторией вычислительных систем Института физики полупроводников им. А.В. Ржанова СО РАН. В библиотеке используется интерфейс профилирования стандарта MPI для перехвата обращений к функциям коллективных операций информационных обменов, остальные вызовы передаются системной библиотеке MPI (MPICH2, OpenMPI и др.). Построение отображений номеров ветвей осуществляется при создании MPI-коммутираторов. Выбор алгоритма реализации ТЦО (функции MPI_Allgather) осуществляется на основе значения переменной среды окружения. Информационные графы алгоритмов формируются в виде списков смежных вершин в упакованном формате CSR – Compressed Sparse Row Graph. Разбиение графа на h подмножеств осуществляется многоуровневым алгоритмом рекурсивной бисекции [7, 8], вычислительная сложность которого составляет $O((|E| + n) \log_2 h)$.

Исследование созданных алгоритмов проводилось на кластерных ВС ЦПВТ ГОУ ВПО «СибГУТИ» и ИВЦ ГОУ ВПО «НГУ». Кластер ЦПВТ ГОУ ВПО «СибГУТИ» построен на базе 10 двухпроцессорных узлов Intel SR2520SAFR с четырехъядерными процессорами Intel Xeon E5420. Первый уровень коммуникационной среды кластера – сеть связи стандарта Gigabit Ethernet, второй уровень – общая память вычислительных узлов. Кластер ИВЦ ГОУ ВПО «НГУ» укомплек-

тован 64 двухпроцессорными узлами Hewlett-Packard BL460c с четырехъядерными процессорами Intel Xeon 5355. Первый уровень коммуникационной среды кластера – сеть связи стандарта InfiniBand 4x DDR, второй уровень – общая память вычислительных узлов.

На кластере ЦИВТ ГОУ ВПО «СибГУТИ» пакет ToroMPI компилировался с библиотекой стандарта MPI – MPICH2 1.2.1 в операционной системе CentOS 5.2 x86-64, а на кластере ИВЦ ГОУ ВПО «НГУ» – в операционной системе SUSE Linux Enterprise Server 10 SP1 x86-64 с библиотекой OpenMPI 1.2.5.

На рис. 4 показано время формирования, разбиений информационных графов алгоритмов и создание отображений $\pi(i)$, $\pi^{-1}(i)$ на процессоре Intel Xeon E5420. Видно, что время, затрачиваемое на формирование отображений номеров ветвей, незначительно. Время же разбиения информационного графа алгоритма Дж. Брука с количеством вершин $n = 1048576$ на $h = 131072$ подмножеств составляет 32,73 с. Поэтому для большемасштабных ВС востребованы параллельные алгоритмы формирования и разбиения графов и создания отображений.

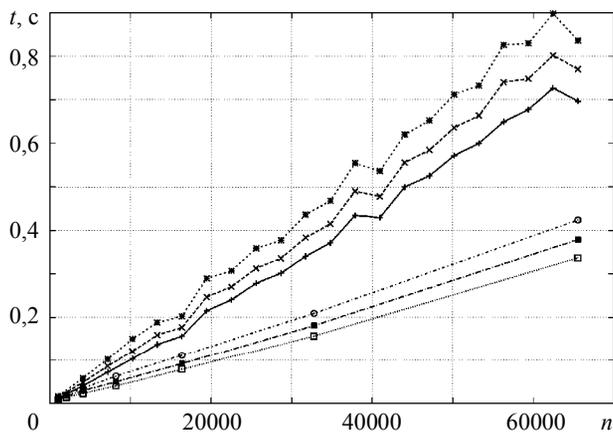


Рис. 4. Зависимость времени разбиения информационных графов на h подмножеств от количества n вершин: —+— — граф алгоритма Дж. Брука, $h = 128$; ----*---- — граф алгоритма Дж. Брука, $h = 256$;*..... — граф алгоритма Дж. Брука, $h = 512$;□..... — граф алгоритма рекурсивного сдваивания, $h = 128$; ---■--- — граф алгоритма рекурсивного сдваивания, $h = 256$;○..... — граф алгоритма рекурсивного сдваивания, $h = 512$

На рис. 5 представлено среднее время ТЦО «короткими» сообщениями на подсистеме из 64 процессорных ядер (4 вычислительных узла). В качестве времени реализации ТЦО бралось максимальное из времен выполнения ветвей параллельной программы. Среднее время рассчитывалось по результатам 100 запусков алгоритма для каждого размера сообщения. Среднее время построения отображений $\pi(i)$ и $\pi^{-1}(i)$ для примера на рис. 5 составило 480 мкс и 300 мкс соответственно для алгоритмов Bruck exch., Bruck reorder и recursive doubling exch., recursive doubling reorder. Видно, что накладные расходы на построение отображений не значительны и компенсируются сокращением времени реализации ТЦО созданными алгоритмами.

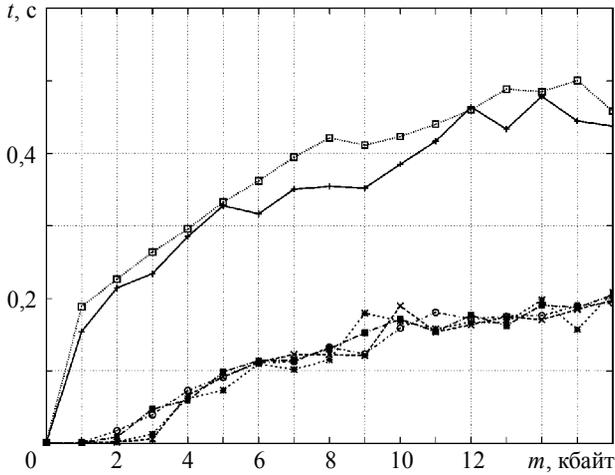


Рис. 5. Зависимость времени ТЦО от размера m сообщения на кластерной ВС ЦПВТ ГОУ ВПО «СибГУТИ» ($n = 64$; $h = 8$; $s_1 = s_2 = \dots = s_8 = 8$): —●— алгоритм Дж. Брука; ----×---- алгоритм Bruck exch.;■..... алгоритм recursive doubling; -.-.-▲-.-.- алгоритм recursive doubling exch.; -.-.-◇-.-.- алгоритма recursive doubling reorder

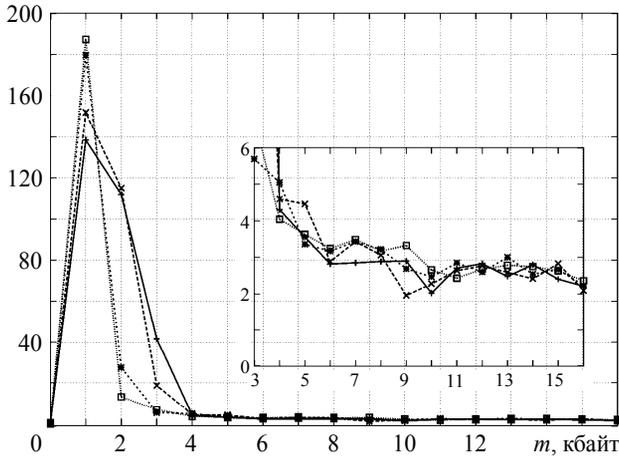


Рис. 6. Зависимость коэффициента ускорения алгоритмов ТЦО от размера m передаваемого сообщения на кластерной ВС ЦПВТ ГОУ ВПО «СибГУТИ» ($n = 64$; $h = 8$; $s_1 = s_2 = \dots = s_8 = 8$): a – ускорение алгоритмов на коротких сообщениях; b – ускорение алгоритмов на длинных сообщениях; —●— алгоритм Bruck exch.; ----×---- алгоритм Bruck reorder;■..... алгоритм recursive doubling exch.; -.-.-▲-.-.- алгоритма recursive doubling reorder

На рис. 6 приведены значения коэффициентов ускорения созданных алгоритмов при передаче коротких и длинных сообщений относительно алгоритмов рекурсивного сдваивания и Дж. Брука. Ускорение в 130 – 180 раз на коротких сообщениях объясняется тем, что при реализации ТЦО сообщениями размером

1 кбайт осуществляется обмен блоками 1, 2, 4, 8, 16 и 32 кбайт, время передачи которых через сеть связи Gigabit Ethernet и общую память узла отличается в 10 – 55 раз (рис. 1, а).

Ускорение алгоритмов реализации ТЦО между 64 параллельными ветвями на кластерной ВС ИВЦ ГОУ ВПО «НГУ» (рис. 7) наблюдается только для сообщений размеров 256 байт – 128 кбайт. Это объясняется тем, что при передаче сообщений размеров больше 2^{22} байт сеть InfiniBand (рис. 1, б) превосходит общую память узла по времени передачи данных.

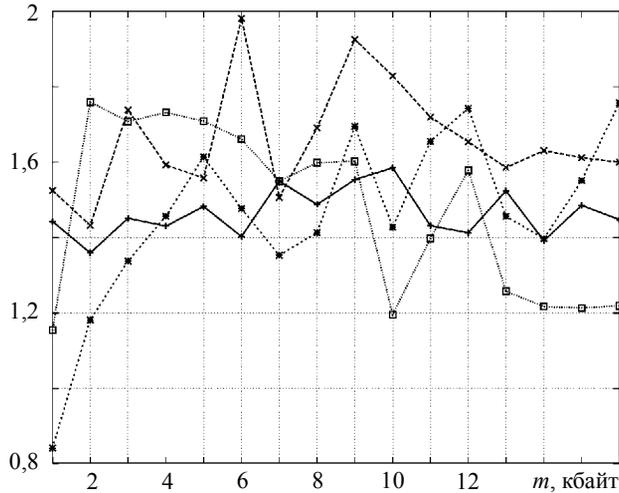


Рис. 7. Ускорение алгоритмов на кластерной ВС ИВЦ ГОУ ВПО «НГУ» ($n = 64$; $h = 8$; $s_1 = s_2 = \dots = s_8 = 8$): а – время ТЦО короткими сообщениями; б – коэффициент ускорение алгоритмов на коротких сообщениях; —◆— – алгоритм Дж. Брука; ----×---- – алгоритм Bruck exch.;■..... – алгоритм Bruck reorder;◆..... – алгоритм рекурсивного сдваивания; ----■---- – алгоритм recursive doubling exch.;◆..... – алгоритма recursive doubling reorder

Созданные алгоритмы целесообразно применять для реализации ТЦО сообщений таких размеров, для которых время их передачи через каналы связи вышележащих уровней превосходит время обменов через каналы связи нижележащих уровней.

Заключение

Предложенный метод позволяет оптимизировать алгоритмы коллективных обменов информацией между ветвями параллельных программ при их реализации на распределенных вычислительных системах с иерархической структурой. Созданные версии алгоритмов рекурсивного сдваивания (Recursive Doubling) и Дж. Брука (J. Bruck) обеспечивают реализацию трансляционно-циклических обменов (ТЦО) с заданным распределением ветвей по элементарным машинам распределенной ВС. Эксперименты на действующих вычислительных кластерах показали превосходство в среднем в 1,1 – 4 раза во времени реализации ТЦО разработанными алгоритмами рекурсивного сдваивания и Дж. Брука по сравнению с

исходными алгоритмами. Накладные расходы на оптимизацию незначительны и компенсируются сокращением времени при многократном обращении к функциям реализации ТЦО.

ЛИТЕРАТУРА

1. Евреинов Э.В., Хорошевский В.Г. Однородные вычислительные систем. Новосибирск: Наука, 1978. 320 с.
2. Хорошевский В.Г. Архитектура вычислительных систем. М.: МГТУ им. Н.Э. Баумана. 2008. 520 с.
3. Rabenseifner R. Automatic MPI counter profiling // Proc. of the 42nd Cray User Group. Noorwijk, 2000. 19 pp.
4. Thakur R., Rabenseifner R., and Gropp W. Optimization of collective communication operations in MPICH // Int. Journal of High Performance Computing Applications. 2005. V. 19. No. 1. P. 49–66.
5. Balaji P. MPI on a million processors / P. Balaji et. al. // Proc. of the PVM/MPI. Berlin: Springer-Verlag, 2009. P. 20–30.
6. Хорошевский В.Г., Курносков М.Г. Алгоритмы распределения ветвей параллельных программ по процессорным ядрам вычислительных систем // Автометрия. 2008. Т. 44. № 2. С. 56–67.
7. Khoroshevsky V., Kurnosov M. Mapping parallel programs into hierarchical distributed computer systems // Proc. of «Software and Data Technologies». Sofia: INSTICC. 2009. V. 2. P. 123–128.
8. Karypis G. and Kumar V. A fast and highly quality multilevel scheme for partitioning irregular graphs // SIAM Journal on Scientific Computing. 1999. V. 20. No. 1. P. 359–392.

Курносков Михаил Георгиевич
ГОУ ВПО «Сибирский государственный университет
телекоммуникаций и информатики»
E-mail: mkurnosov@gmail.com

Поступила в редакцию 3 декабря 2010 г.