

УДК 519.6

*Н. А. ЛИХОДЕД, М. А. ПОЛЕЩУК***МЕТОД РАНЖИРОВАНИЯ ПАРАМЕТРОВ РАЗМЕРА БЛОКОВ ВЫЧИСЛЕНИЙ
ПАРАЛЛЕЛЬНОГО АЛГОРИТМА***(Представлено членом-корреспондентом Л. А. Яновичем)**Белорусский государственный университет, Минск, Беларусь**likhoded@bsu.by; poleschuma@bsu.by*

Исследуется задача получения макроопераций параллельного алгоритма, приводящих к меньшему числу обращений к глобальной памяти. Сформулированы и доказаны утверждения, позволяющие оценить объем коммуникационных операций, порождаемых разбиением множества итераций.

Ключевые слова: параллельные вычисления, распараллеливание алгоритмов, графический процессор, минимизация объема коммуникационных операций.

*N. A. LIKHODED, M. A. PALIASHCHUK***METHOD OF RANKING TILES SIZE PARAMETERS OF A PARALLEL ALGORITHM***Belarusian State University, Minsk, Belarus**likhoded@bsu.by; poleschuma@bsu.by*

A method for obtaining tiles of operations of a parallel algorithm is developed. Propositions for ranking tiles size parameters are stated and proved. Statements to assess the amount of communication operations generated by the partition of the set of iterations are stated and proved.

Keywords: parallel computing, parallelization of algorithms, GPU, minimization of communications.

Введение. Время решения задачи на современном компьютере во многом определяется степенью использования памяти с быстрым доступом. В качестве целевого компьютера, т. е. компьютера, на котором требуется реализовать параллельную версию алгоритма, будем рассматривать графические процессоры (GPU). При многопроцессорной обработке на GPU быстрым является процесс обращения к разделяемой памяти мультипроцессора, оперирующего некоторыми данными алгоритма, но не обращение к глобальной памяти GPU. Чем меньше число обращений к глобальной памяти, тем быстрее выполняется алгоритм.

Для реализации алгоритма на графическом процессоре множество операций алгоритма должно быть разбито на блоки, а блоки – на потоки (нити) вычислений. Множество операций алгоритма разбить на блоки можно путем тайлинга (тайлинг первого уровня), затем разбить блоки вычислений на потоки вычислений можно путем повторного применения тайлинга (тайлинг второго уровня) [1]. Тайлинг (tiling) – это преобразование алгоритма для получения макроопераций-тайлов [2; 3]; операции одного тайла выполняются атомарно, как одна единица вычислений, а обмен данными происходит массивами.

В данной работе предлагается метод, позволяющий получать тайлы вычислений с меньшим числом обращений к глобальной памяти. Используется анализ информационных зависимостей, порождающих коммуникационные операции. Сформулированы и доказаны утверждения, позволяющие ранжировать параметры размера тайлов. Предлагаемый метод можно также применять после некоторых предварительных оптимизирующих преобразований алгоритма. В частности, ранжировать параметры размера тайлов предлагаемым способом можно после преобразо-

ваний для улучшения параллелизма вычислений уровня тайлов и обеспечения одновременного начала шаблонных вычислений процессорами [4; 5]. Результаты исследований этой работы могут быть использованы для минимизации объема коммуникационных операций между параллельными вычислительными процессами, реализуемыми на суперкомпьютерах с распределенной памятью.

Предварительные сведения. Приведем необходимые для дальнейшего изложения сведения о формальном описании алгоритма и о тайлинге.

Будем считать, что алгоритм задан последовательной программой линейного класса [6]. Основную вычислительную часть такой программы составляют циклические конструкции; границы изменения параметров циклов задаются неоднородными формами, линейными по совокупности параметров циклов и внешних переменных. Предполагается, что в гнезде циклов имеется K выполняемых операторов S_β и используется L массивов a_l размерностей v_l . Область изменения параметров циклов (область итераций) для оператора S_β и размерность этой области обозначим соответственно V_β и n_β .

Выполнение оператора S_β при конкретных значениях β вектора параметров цикла J будем называть операцией и обозначать $S_\beta(J)$. Зависимости (информационные связи) между операциями задаются функциями вида

$$\bar{\Phi}_{\alpha,\beta}(J) = \Phi_{\alpha,\beta}J + \Psi_{\alpha,\beta}N - \varphi^{\alpha,\beta}, \quad (1)$$

$$J \in V_{\alpha,\beta}, N \in \mathbb{Z}^s, \Phi_{\alpha,\beta} \in \mathbb{Z}^{n_\alpha \times n_\beta}, \Psi_{\alpha,\beta} \in \mathbb{Z}^{n_\alpha \times s}, \varphi^{\alpha,\beta} \in \mathbb{Z}^{n_\alpha},$$

где $N \in \mathbb{Z}^s$ – вектор внешних переменных алгоритма; s – число внешних переменных. Функция зависимостей $\bar{\Phi}_{\alpha,\beta}(J)$ позволяет для операции $S_\beta(J)$ найти операцию $S_\alpha(I)$, от которой $S_\beta(J)$ зависит. Функции зависимостей являются удобным математическим аппаратом для описания информационных связей между операциями алгоритма (другие названия этих функций: покрывающие функции графа алгоритма [6], h -преобразования [7; 8]).

Вхождением (l, β, q) будем называть q -е вхождение массива a_l в оператор S_β . Индексы элементов l -го массива, связанных с вхождением (l, β, q) , выражаются функцией $\bar{F}_{l,\beta,q}$ вида

$$\bar{F}_{l,\beta,q}(J) = F_{l,\beta,q}J + G_{l,\beta,q}N + f^{l,\beta,q},$$

$$J(j_1, \dots, j_{n_\beta}) \in V_\beta, N \in \mathbb{Z}^s, F_{l,\beta,q} \in \mathbb{Z}^{v_l \times n_\beta}, G_{l,\beta,q} \in \mathbb{Z}^{v_l \times s}, f^{l,\beta,q} \in \mathbb{Z}^{v_l}.$$

Пара вхождений $(l, \alpha, 1)$ и (l, β, q) порождает истинную зависимость $S_\alpha(I) \rightarrow S_\beta(J)$, если $S_\alpha(I)$ выполняется раньше $S_\beta(J)$; $S_\alpha(I)$ переопределяет (изменяет) элемент массива a_l , а $S_\beta(J)$ использует в качестве аргумента тот же элемент массива; между операциями $S_\alpha(I)$ и $S_\beta(J)$ этот элемент не переопределяется.

Пусть в гнезде циклов имеется Θ наборов выполняемых операторов. Под набором операторов будем понимать один или несколько операторов, окруженных одним и тем же множеством циклов. Операторы и наборы операторов линейно упорядочены расположением их в записи алгоритма. Обозначим: $V^\vartheta, 1 \leq \vartheta \leq \Theta$, – области изменения параметров циклов, окружающих наборы операторов, n^ϑ – размерность области V^ϑ , число циклов, окружающих ϑ -й набор операторов. Заметим, что если оператор S_β , принадлежит набору операторов с номером ϑ^β , то область V_β может быть уже области V^{ϑ^β} .

Как уже отмечалось, тайлинг – это преобразование алгоритма для получения макроопераций, называемых зерном вычислений, или тайлами. При тайлинге каждый цикл разбивается на два цикла: глобальный, параметр которого определяет на данном уровне вложенности порядок вычисления тайлов, и локальный, в котором параметр исходного цикла изменяется в границах одного тайла. Допускается вырожденное разбиение цикла, при котором все итерации относятся к глобальному циклу или все итерации относятся к локальному циклу.

Следующие величины и множества используются для формализации тайлинга.

$m_\zeta^\vartheta = \min_{J(j_1, j_2, \dots, j_{n^\vartheta}) \in V^\vartheta} j_\zeta, M_\zeta^\vartheta = \max_{J(j_1, j_2, \dots, j_{n^\vartheta}) \in V^\vartheta} j_\zeta, 1 \leq \zeta \leq n^\vartheta$, – предельные значения изменения параметров циклов;

$r_1^{\vartheta}, \dots, r_{n^{\vartheta}}^{\vartheta}$ – заданные натуральные числа, определяющие размеры тайла; r_{ζ}^{ϑ} обозначает число значений параметра j_{ζ} , приходящихся на один тайл ϑ -го набора операторов; r_{ζ}^{ϑ} может принимать фиксированное значение в пределах от 1 до $r_{\zeta}^{\vartheta, \max}$ включительно, где $r_{\zeta}^{\vartheta, \max} = M_{\zeta}^{\vartheta} - m_{\zeta}^{\vartheta} + 1$; если $r_{\zeta}^{\vartheta} = 1$, то цикл с параметром j_{ζ} является глобальным не разбиваемым; если $r_{\zeta}^{\vartheta} = r_{\zeta}^{\vartheta, \max}$, то цикл с параметром j_{ζ} является локальным не разбиваемым; если два набора операторов имеют общий цикл с параметром j_{ζ} , то $r_{\zeta}^{\vartheta 1} = r_{\zeta}^{\vartheta 2}$;

$Q_{\zeta}^{\vartheta} = \lceil (M_{\zeta}^{\vartheta} - m_{\zeta}^{\vartheta} + 1) / r_{\zeta}^{\vartheta} \rceil$, $1 \leq \zeta \leq n^{\vartheta}$, – число частей, на которые при формировании тайлов разбивается область значений параметра j_{ζ} цикла, окружающего ϑ -й набор операторов;

$V^{\vartheta, \text{gl}} = \{J^{\text{gl}}(j_1^{\text{gl}}, \dots, j_{n^{\vartheta}}^{\text{gl}}) \mid 0 \leq j_{\zeta}^{\text{gl}} \leq Q_{\zeta}^{\vartheta} - 1, 1 \leq \zeta \leq n^{\vartheta}\}$ – области изменения параметров глобальных, т. е. уровня тайлов, циклов;

$V_{J^{\text{gl}}}^{\vartheta} = \{J(j_1, \dots, j_{n^{\vartheta}}) \in V^{\vartheta} \mid m_{\zeta}^{\vartheta} + j_{\zeta}^{\text{gl}} r_{\zeta}^{\vartheta} \leq j_{\zeta} \leq m_{\zeta}^{\vartheta} - 1 + (j_{\zeta}^{\text{gl}} + 1) r_{\zeta}^{\vartheta}, 1 \leq \zeta \leq n^{\vartheta}\}$, $J^{\text{gl}} \in V^{\vartheta, \text{gl}}$, – области изменения параметров локальных (уровня операций тайлов) циклов при фиксированных значениях параметров глобальных циклов. Множество операций, выполняемых на итерациях множества $V_{J^{\text{gl}}}^{\vartheta}$, будем также обозначать $V_{J^{\text{gl}}}^{\vartheta}$. Множества $V_{J^{\text{gl}}}^{\vartheta}$ называются тайлами.

Условия, характеризующие объем коммуникационных операций. Результаты, представленные в этом разделе, позволяют оценить объем коммуникационных операций между тайлами первого уровня (т. е. между блоками вычислений).

Пусть рассматривается гнездо циклов, для которого тайлинг допустим. Если заведомо известно, что $r_{\zeta}^{\vartheta} = 1$ (цикл с параметром j_{ζ} является глобальным не разбиваемым) или $r_{\zeta}^{\vartheta} = r_{\zeta}^{\vartheta, \max}$ (цикл с параметром j_{ζ} является локальным не разбиваемым), то разбиение итераций j_{ζ} не рассматривается. Обозначим через $(\Phi_{\alpha, \beta})_{\zeta}$ и $(\Psi_{\alpha, \beta})_{\zeta}$ строки матриц $\Phi_{\alpha, \beta}$ и $\Psi_{\alpha, \beta}$ с номером ζ , через $e_{\zeta}^{(n_{\beta})}$ – вектор-строку размера n_{β} , у которой координата с номером ζ равна 1, а остальные координаты нулевые; обозначим еще $\eta_{\zeta}^{\alpha, \beta} = (\Psi_{\alpha, \beta})_{\zeta} N - \Phi_{\zeta}^{\alpha, \beta}$.

Л е м м а 1. Пусть определение элемента некоторого массива a_l происходит на вхождении $(l, \alpha, 1)$ в левой части оператора S_{α} , а использование – на вхождении (l, β, q) в правой части оператора S_{β} , причем в окружении операторов S_{α} и S_{β} имеется цикл с параметром j_{ζ} , и выполняется условие

$$(\Phi_{\alpha, \beta})_{\zeta} = e_{\zeta}^{(n_{\beta})}. \quad (2)$$

Если

$$\eta_{\zeta}^{\alpha, \beta} = 0, \quad (3)$$

то определение и использование элемента массива данных происходит при одном и том же значении параметра цикла j_{ζ}^{gl} .

Если

$$0 < |\eta_{\zeta}^{\alpha, \beta}| < r_{\zeta}^{\vartheta}, \quad (4)$$

то определение и использование элемента массива данных происходит при одном и том же значении параметра цикла j_{ζ}^{gl} на $r_{\zeta}^{\vartheta} - |\eta_{\zeta}^{\alpha, \beta}|$ итерациях из каждых r_{ζ}^{ϑ} итераций j_{ζ} , а на $|\eta_{\zeta}^{\alpha, \beta}|$ итерациях из каждых r_{ζ}^{ϑ} итераций j_{ζ} определение и использование элемента массива данных происходит при разных значениях параметра цикла j_{ζ}^{gl} .

Если

$$|\eta_{\zeta}^{\alpha, \beta}| \geq r_{\zeta}^{\vartheta}, \quad (5)$$

то определение и использование элемента массива данных происходит при разных значениях параметра цикла j_{ζ}^{gl} .

Д о к а з а т е л ь с т в о. Докажем первое утверждение леммы. Функция зависимостей $I = \Phi_{\alpha, \beta}(J)$ определяет операцию $S_{\alpha}(I(i_1, \dots, i_{n_{\alpha}}))$ вычисления данного, требуемого операцией $S_{\beta}(J(j_1, \dots, j_{n_{\beta}}))$ в качестве аргумента. Из равенств (1)–(3) следует $i_{\zeta} = (\Phi_{\alpha, \beta})_{\zeta} J = e_{\zeta}^{(n_{\beta})} J = j_{\zeta}$, поэтому значения параметров цикла j_{ζ}^{gl} , на которых определяется и используется элемент массива данных, совпадают.

Докажем второе утверждение леммы. Так как из определения тайлов имеем

$$j_\zeta = m_\zeta^{\mathfrak{g}} + j_\zeta^{\mathfrak{gl}} r_\zeta^{\mathfrak{g}} + \alpha_\zeta^{\mathfrak{g}},$$

где $0 \leq \alpha_\zeta^{\mathfrak{g}} < r_\zeta^{\mathfrak{g}}$,
то

$$i_\zeta = (\Phi_{\alpha,\beta})_\zeta J + \eta_\zeta^{\alpha,\beta} = e_\zeta^{(n\beta)} J + \eta_\zeta^{\alpha,\beta} = m_\zeta^{\mathfrak{g}} + j_\zeta^{\mathfrak{gl}} r_\zeta^{\mathfrak{g}} + \alpha_\zeta^{\mathfrak{g}} + \eta_\zeta^{\alpha,\beta}.$$

Найдем значение параметра цикла $i_\zeta^{\mathfrak{gl}}$, на котором определяется элемент массива данных ($\lfloor \cdot \rfloor$ обозначает ближайшее «снизу» целое число):

$$i_\zeta^{\mathfrak{gl}} = \left\lfloor \frac{i_\zeta - m_\zeta^{\mathfrak{g}}}{r_\zeta^{\mathfrak{g}}} \right\rfloor = \left\lfloor \frac{j_\zeta^{\mathfrak{gl}} r_\zeta^{\mathfrak{g}} + \alpha_\zeta^{\mathfrak{g}} + \eta_\zeta^{\alpha,\beta}}{r_\zeta^{\mathfrak{g}}} \right\rfloor = j_\zeta^{\mathfrak{gl}} + \left\lfloor \frac{\alpha_\zeta^{\mathfrak{g}} + \eta_\zeta^{\alpha,\beta}}{r_\zeta^{\mathfrak{g}}} \right\rfloor. \quad (6)$$

Последний член равен нулю при $r_\zeta^{\mathfrak{g}} - |\eta_\zeta^{\alpha,\beta}|$ значениях $\alpha_\zeta^{\mathfrak{g}}$, на остальных $|\eta_\zeta^{\alpha,\beta}|$ значениях $\alpha_\zeta^{\mathfrak{g}}$ последний член не равен нулю.

Для доказательства третьего утверждения теоремы достаточно заметить, что если выполняется условие (5), то последний член в равенстве (6) не является нулевым. \square

В частном случае $\Psi = 0$ первые два утверждения теоремы 1 сформулированы в работе [9].

Л е м м а 2. Пусть определение элемента некоторого массива a_l происходит при выполнении операции $S_\alpha(I(i_1, \dots, i_{n_\alpha}))$, а использование – при выполнении операции $S_\beta(J(j_1, \dots, j_{n_\beta}))$, причем в окружении операторов S_α и S_β имеется цикл с параметром j_ζ . Если выполняется условие

$$(\Phi_{\alpha,\beta})_\zeta = \alpha_\zeta e_\zeta^{(n\beta)}, \quad (7)$$

где $\alpha_\zeta \in \mathbb{Z}$, $\alpha_\zeta \neq 1$,

то для всех итераций j_ζ , удовлетворяющих условию

$$|(\alpha_\zeta - 1)j_\zeta + \eta_\zeta^{\alpha,\beta}| \geq r_\zeta^{\mathfrak{g}}, \quad (8)$$

определение и использование элемента массива данных происходит при разных значениях параметра цикла $j_\zeta^{\mathfrak{gl}}$.

Д о к а з а т е л ь с т в о. Преобразуем выражение под знаком модуля в неравенстве (8):

$$\begin{aligned} (\alpha_\zeta - 1)j_\zeta + \eta_\zeta^{\alpha,\beta} &= \alpha_\zeta j_\zeta + \eta_\zeta^{\alpha,\beta} - j_\zeta = \alpha_\zeta e_\zeta^{(n\beta)} J + \eta_\zeta^{\alpha,\beta} - j_\zeta = \\ (\Phi_{\alpha,\beta})_\zeta J + \eta_\zeta^{\alpha,\beta} - j_\zeta &= (\Phi_{\alpha,\beta})_\zeta J + (\Psi_{\alpha,\beta})_\zeta N - \varphi_\zeta^{\alpha,\beta} - j_\zeta = i_\zeta - j_\zeta. \end{aligned}$$

Таким образом, если $\alpha_\zeta \neq 1$, то $|i_\zeta - j_\zeta| \geq r_\zeta^{\mathfrak{g}}$. Поэтому определение и использование элемента массива данных происходит при разных значениях параметра цикла $j_\zeta^{\mathfrak{gl}}$. \square

З а м е ч а н и е 1. Из леммы 2 следует, что при выполнении условия (7) число итераций j_ζ , для которых определение и использование элемента массива данных происходит при одном и том же значении цикла с параметром $j_\zeta^{\mathfrak{gl}}$, не может быть большим (число $r_\zeta^{\mathfrak{g}}$, как правило, гораздо меньше количества всех итераций j_ζ). Основная причина этого заключается в том, что модуль разности $|j_\zeta - i_\zeta|$ не ограничен небольшим (в пределах нескольких единиц) числом. В случае выполнения условия

$$(\Phi_{\alpha,\beta})_\zeta \neq \alpha_\zeta e_\zeta^{(n\beta)},$$

где $\alpha_\zeta \in \mathbb{Z}$,

модуль разности $|j_\zeta - i_\zeta|$ также не ограничен небольшим числом, так как зависит от J . Поэтому лишь для небольшого числа итераций j_ζ определение и использование элемента массива данных может происходить при одинаковых значениях цикла с параметром $j_\zeta^{\mathfrak{gl}}$.

Пусть вхождение (l, β, q) в правую часть некоторого оператора порождает истинную зависимость, $\overline{\Phi}_{\alpha,\beta}$ – функция зависимостей. Обозначим

$$\rho_{l,\beta,q} = \text{rank } F_{l,\beta,q}, \rho_{l,\beta,q}^\zeta = \text{rank} \begin{pmatrix} F_{l,\beta,q} \\ e_\zeta^{(n\beta)} \end{pmatrix}, \rho_{l,\beta,q}^\Phi = \text{rank} \begin{pmatrix} F_{l,\beta,q} \\ \Phi_{\alpha,\beta} \end{pmatrix}, \rho_{l,\beta,q}^{\Phi,\zeta} = \text{rank} \begin{pmatrix} F_{l,\beta,q} \\ \Phi_{\alpha,\beta} \\ e_\zeta^{(n\beta)} \end{pmatrix}.$$

Если вхождение (l, β, q) не порождает истинную зависимость, то по определению $\rho_{l,\beta,q}^\Phi = \rho_{l,\beta,q}$, $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\zeta$. Отметим, что $\rho_{l,\beta,q}^\Phi$ и $\rho_{l,\beta,q}^{\Phi,\zeta} - 1$ отличаются не более чем на 1.

Обозначим $M^\vartheta = \max_{\zeta} (M_\zeta^\vartheta - m_\zeta^\vartheta) + 1$ – наибольшее число итераций циклов, участвующих в получении тайлов. Для простоты записи будем использовать обозначение M без индекса ϑ , где будет подразумеваться набор операторов \mathfrak{G}^β при упоминании оператора S_β . Отметим, что величина $Q_\zeta^\vartheta r_\zeta^\vartheta$ имеет порядок M . Определим термин «фиксированное данное массива» как конкретное, неизмененное содержимое соответствующей ячейки памяти. Следующая лемма позволяет оценить число используемых фиксированных данных на каждом вхождении (l, β, q) при фиксированных значениях циклов с параметрами j_ζ и j_ζ^{gl} , общее число используемых фиксированных данных на вхождении (l, β, q) .

Л е м м а 3. Число фиксированных данных, используемых на вхождении (l, β, q) в правой части оператора S_β , при фиксированном значении цикла с параметром j_ζ оценивается величиной $O(M^{\rho_{l,\beta,q}^{\Phi,\zeta}-1})$.

Число фиксированных данных, используемых на вхождении (l, β, q) , оценивается величиной $O(M^{\rho_{l,\beta,q}^\Phi})$.

Число фиксированных данных, используемых на вхождении (l, β, q) , при фиксированном значении глобального цикла с параметром j_ζ^{gl} оценивается величиной $O(M^{\rho_{l,\beta,q}^\Phi})$, если $\rho_{l,\beta,q}^{\Phi,\zeta} - 1 = \rho_{l,\beta,q}^\Phi$, и величиной $O(r_\zeta^\vartheta M^{\rho_{l,\beta,q}^\Phi - 1})$, если $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\Phi$.

Д о к а з а т е л ь с т в о. На вхождении (l, β, q) фиксированное данное используется при фиксированном значении цикла с параметром j_ζ на итерациях подпространства итераций размерности $n_\beta - \rho_{l,\beta,q}^{\Phi,\zeta}$ [10; 11]. Тогда число фиксированных данных, используемых на вхождении (l, β, q) при фиксированном значении цикла с параметром j_ζ , т. е. число подпространств итераций размерности $n_\beta - \rho_{l,\beta,q}^{\Phi,\zeta}$ в пространстве размерности $n_\beta - 1$ есть $O(M^{n_\beta - 1 - (n_\beta - \rho_{l,\beta,q}^{\Phi,\zeta})}) = O(M^{\rho_{l,\beta,q}^{\Phi,\zeta} - 1})$. Аналогично, число фиксированных данных, используемых на вхождении (l, β, q) , есть число подпространств итераций размерности $n_\beta - \rho_{l,\beta,q}^\Phi$ в пространстве размерности n_β и для них справедлива асимптотическая оценка $O(M^{n_\beta - (n_\beta - \rho_{l,\beta,q}^\Phi)}) = O(M^{\rho_{l,\beta,q}^\Phi})$.

Если выполняется условие $\rho_{l,\beta,q}^{\Phi,\zeta} - 1 = \rho_{l,\beta,q}^\Phi$, то оценка числа используемых на вхождении (l, β, q) фиксированных данных одинакова как для одного j_ζ , так и для любого числа итераций j_ζ . Поэтому число фиксированных данных при r_ζ^ϑ значениях цикла с параметром j_ζ оценивается величиной $O(M^{\rho_{l,\beta,q}^{\Phi,\zeta} - 1}) = O(M^{\rho_{l,\beta,q}^\Phi})$. Если выполняется условие $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\Phi$, то для каждого j_ζ (всего r_ζ^ϑ значений j_ζ) используются новые фиксированные данные. Так как $Q_\zeta^\vartheta r_\zeta^\vartheta = O(M)$, то r_ζ^ϑ может быть величиной порядка M и ее следует учесть в асимптотической оценке $O(r_\zeta^\vartheta M^{\rho_{l,\beta,q}^{\Phi,\zeta} - 1})$ числа фиксированных данных. \square

Оценка объема коммуникационных операций. В этом разделе для каждого вхождения (l, β, q) в правую часть оператора S_β оценивается объем коммуникационных операций чтения и соответствующих им операций записи, порождаемых разбиением итераций j_ζ при получении блоков вычислений.

Сделаем предположения, необходимые для практического использования результатов исследований: элементы матриц $\Phi_{\alpha,\beta}$, $\Psi_{\alpha,\beta}$ и векторов $\varphi^{\alpha,\beta}$ по модулю не превосходят нескольких единиц; $|\eta_\zeta^{\alpha,\beta}| = |\varphi_\zeta^{\alpha,\beta}| < r_\zeta^\vartheta$, если $(\Psi_{\alpha,\beta})_\zeta = 0$; $|\eta_\zeta^{\alpha,\beta}| \geq r_\zeta^\vartheta$, если $(\Psi_{\alpha,\beta})_\zeta \neq 0$.

Т е о р е м а. Пусть вхождение (l, β, q) порождает истинную зависимость: определение некоторого данного происходит на вхождении $(l, \alpha, 1)$ в левой части оператора S_α , а использование – на вхождении (l, β, q) в правой части оператора S_β , причем в окружении операторов S_α и S_β имеется цикл с параметром j_ζ .

Тогда при получении блоков вычислений для реализации алгоритма на графическом процессоре, разбиение итераций цикла j_ζ порождает коммуникационные операции чтения и записи, объем которых имеет следующие оценки:

1) если выполняются условия (2), (3), то не требуется ни операций чтения, ни операций записи;

2) если выполняются условия (2), (4), то требуется $O(Q_\zeta^g M^{\rho_{l,\beta,q}^\Phi}^{-1})$ операций чтения и операций записи;

3) если условие (2) не выполняется, выполняется условие $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\Phi$, то требуется $O(M^{\rho_{l,\beta,q}^\Phi})$ операций чтения и операций записи;

4) если условие (2) не выполняется, выполняется условие $\rho_{l,\beta,q}^{\Phi,\zeta} - 1 = \rho_{l,\beta,q}^\Phi$, то требуется $O(Q_\zeta^g M^{\rho_{l,\beta,q}^\Phi})$ операций чтения и $O(M^{\rho_{l,\beta,q}^\Phi})$ операций записи.

В случае, когда вхождение (l, β, q) не порождает истинной зависимости (происходит обращение к входным данным) или цикл с параметром j_ζ имеется в окружении только оператора S_β , оценки следующие:

5) если выполняется условие $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\Phi$, то требуется $O(M^{\rho_{l,\beta,q}^\Phi})$ операций чтения;

6) если выполняется условие $\rho_{l,\beta,q}^{\Phi,\zeta} - 1 = \rho_{l,\beta,q}^\Phi$, то требуется $O(Q_\zeta^g M^{\rho_{l,\beta,q}^\Phi})$ операций чтения.

Д о к а з а т е л ь с т в о. Утверждения теоремы оценивают объем коммуникационных операций, порождаемых разбиением итераций цикла j_ζ , в зависимости от возможных значений координат строки $(\Phi_{\alpha,\beta})_\zeta$ и величины $\eta_\zeta^{\alpha,\beta}$ (шесть случаев, возможных на практике). Напомним, множество итераций цикла с параметром j_ζ разбивается на Q_ζ^g частей.

Рассмотрим утверждение первого случая теоремы. Пусть вхождение (l, β, q) порождает истинную зависимость и выполняются условия (2), (3). Тогда выполняется равенство $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\Phi$; из леммы 1 и леммы 3 следует, что в каждой из Q_ζ^g частей все $O(r_\zeta^g M^{\rho_{l,\beta,q}^\Phi}^{-1})$ фиксированных данных определяются и используются при одном и том же значении параметра цикла j_ζ^{gl} . Коммуникационных операций не требуется.

Во втором случае теоремы вхождение (l, β, q) порождает истинную зависимость, выполняются условия (2), (4). Тогда выполняется равенство $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\Phi$; из леммы 1 и леммы 3 следует, что в каждой из Q_ζ^g частей $|\eta_\zeta^{\alpha,\beta}| O(M^{\rho_{l,\beta,q}^\Phi}^{-1})$ данных, из общего числа $O(r_\zeta^g M^{\rho_{l,\beta,q}^\Phi}^{-1})$ используемых в каждой части фиксированных данных, определяются и используются при разных итерациях цикла j_ζ^{gl} , причем независимо от значения r_ζ^g величина $|\eta_\zeta^{\alpha,\beta}|$ не превышает нескольких единиц. Суммарный объем коммуникационных операций чтения и операций записи по всем Q_ζ^g частям определяется оценкой $Q_\zeta^g O(M^{\rho_{l,\beta,q}^\Phi}^{-1}) = O(Q_\zeta^g M^{\rho_{l,\beta,q}^\Phi}^{-1})$.

Рассмотрим третий и четвертый случаи теоремы. В третьем случае выполняется условие $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\Phi$, в четвертом – условие $\rho_{l,\beta,q}^{\Phi,\zeta} - 1 = \rho_{l,\beta,q}^\Phi$. Согласно лемме 3, число используемых фиксированных данных в каждой из Q_ζ^g частей оценивается в третьем случае величиной $O(r_\zeta^g M^{\rho_{l,\beta,q}^\Phi}^{-1})$, в четвертом – величиной $O(M^{\rho_{l,\beta,q}^\Phi})$. По условию, вхождение (l, β, q) порождает истинную зависимость, но условие (2) не выполняется. Поэтому в третьем случае нельзя утверждать, что в каждой из Q_ζ^g частей только $O(M^{\rho_{l,\beta,q}^\Phi}^{-1})$ данных из общего числа $O(r_\zeta^g M^{\rho_{l,\beta,q}^\Phi}^{-1})$ используемых (в каждой части) фиксированных данных определяются и используются при разных итерациях цикла j_ζ^{gl} (см. замечание 1); объемы коммуникационных операций чтения и операций записи по всем Q_ζ^g частям следует оценить величиной $Q_\zeta^g O(r_\zeta^g M^{\rho_{l,\beta,q}^\Phi}^{-1}) = O(M^{\rho_{l,\beta,q}^\Phi})$. Аналогично, в четвертом случае объем коммуникационных операций чтения по всем Q_ζ^g частям равен $Q_\zeta^g O(M^{\rho_{l,\beta,q}^\Phi}) = O(Q_\zeta^g M^{\rho_{l,\beta,q}^\Phi})$. Объем коммуникационных операций записи в четвертом случае оценивается величиной $O(M^{\rho_{l,\beta,q}^\Phi})$ – числом фиксированных данных, используемых, согласно лемме 3, на вхождении (l, β, q) как на одной части (выполняется условие $\rho_{l,\beta,q}^{\Phi,\zeta} - 1 = \rho_{l,\beta,q}^\Phi$), так и на всех Q_ζ^g частях.

В пятом и шестом случаях теоремы вхождение (l, β, q) не порождает истинной зависимости (происходит обращение к входным данным) или порождает истинную зависимость $S_\alpha(I) \rightarrow S_\beta(J)$, но цикл с параметром j_ζ имеется в окружении только оператора S_β . В пятом случае выполняется условие $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\Phi$, в шестом – условие $\rho_{l,\beta,q}^{\Phi,\zeta} - 1 = \rho_{l,\beta,q}^\Phi$. Число используемых фиксированных

данных в каждой из Q_ζ^9 частей оценивается, согласно лемме 3, в пятом случае величиной $O(r_\zeta^9 M^{\rho_{l,\beta,q}^\Phi})$, в шестом – величиной $O(M^{\rho_{l,\beta,q}^\Phi})$. Тогда объем коммуникационных операций (только чтение) по всем Q_ζ^9 частям в пятом случае равен $Q_\zeta^9 O(r_\zeta^9 M^{\rho_{l,\beta,q}^\Phi}) = O(M^{\rho_{l,\beta,q}^\Phi})$, в шестом – $Q_\zeta^9 O(M^{\rho_{l,\beta,q}^\Phi}) = O(Q_\zeta^9 M^{\rho_{l,\beta,q}^\Phi})$. \square

З а м е ч а н и е 2. Если вхождение (l, β, q) не порождает истинной зависимости, то ранее введены обозначения $\rho_{l,\beta,q}^\Phi = \rho_{l,\beta,q}$ и $\rho_{l,\beta,q}^{\Phi,\zeta} = \rho_{l,\beta,q}^\zeta$. Поэтому предположения пунктов 5 и 6 теоремы можно записать в виде $\rho_{l,\beta,q}^\zeta = \rho_{l,\beta,q}$ и $\rho_{l,\beta,q}^\zeta - 1 = \rho_{l,\beta,q}$, а объем операций чтения в виде $O(M^{\rho_{l,\beta,q}^\zeta})$ и $O(Q_\zeta^9 M^{\rho_{l,\beta,q}^\zeta})$ соответственно.

Приоритеты разбиения итераций циклов. Утверждения теоремы позволяют выяснить асимптотику суммарного объема коммуникационных операций, порождаемых разбиением множества итераций j_ζ на Q_ζ^9 частей, в зависимости от возможных значений координат строки $(\Phi_{\alpha,\beta})_\zeta$ и величины $\eta_\zeta^{\alpha,\beta}$. Полученные оценки позволяют ранжировать параметры размера тайлов для минимизации объема коммуникационных операций между блоками вычислений.

Для каждого вхождения (l, β, q) в правую часть оператора S_β , окруженного циклом с параметром j_ζ , утверждения теоремы определяют объемы коммуникационных операций чтения $\omega_{l,\beta,q}^{\zeta,R}$ и операций записи $\omega_{l,\beta,q}^{\zeta,W}$. Их величины могут принимать одно из значений: нуль (первый случай теоремы); $O(Q_\zeta^9 M^\tau)$, где $\tau \in Z, \tau \geq 0$, (второй, четвертый и шестой случаи теоремы); $O(M^\tau)$ (третий, четвертый и пятый случаи теоремы). Для координаты j_ζ обозначим через ω_ζ^R и ω_ζ^W сумму объемов коммуникационных операций чтения и операций записи по всем вхождениям (l, β, q) : $\omega_\zeta^R = \sum_{(l,\beta,q)} \omega_{l,\beta,q}^{\zeta,R}$, $\omega_\zeta^W = \sum_{(l,\beta,q)} \omega_{l,\beta,q}^{\zeta,W}$. Определим также суммарный объем ω_ζ всех (чтение и запись) коммуникационных операций для координаты j_ζ : $\omega_\zeta = \omega_\zeta^R + \omega_\zeta^W$. Заметим, что каждая из величин ω_ζ^R , ω_ζ^W и ω_ζ равна нулю или оценивается величинами $O(Q_\zeta^9 M^\tau)$, $O(M^\tau)$ (вид оценки зависит от наличия членов с множителем Q_ζ^9).

Определим через z_ζ^R , z_ζ^W и z_ζ приоритеты разбиения цикла с параметром j_ζ по чтению, по записи и суммарно по чтению и записи. Положим z_ζ^R равным -1 , если $\omega_\zeta^R = 0$; равным τ , если $\omega_\zeta^R = O(M^\tau)$; равным $\tau + 0,5$, если $\omega_\zeta^R = O(Q_\zeta^9 M^\tau)$; z_ζ^W и z_ζ определяются аналогично через значения ω_ζ^W и ω_ζ . Число $0,5$ отражает зависимость Q_ζ^9 и M , что следует из равенства $M = O(r_\zeta^9 Q_\zeta^9)$. Приоритет -1 является самым высоким (разбиение координаты j_ζ не приводит к коммуникационным операциям), с ростом значения величин z_ζ^R , z_ζ^W и z_ζ приоритет убывает. Отметим, что если в реальных вычислениях на графическом процессоре скорости доступа к глобальной памяти по чтению и по записи (или их средние оценки) известны и они различаются на один или более порядок относительно M , следует вместо суммарного по чтению и записи приоритета z_ζ использовать приоритеты по чтению z_ζ^R и по записи z_ζ^W с учетом такого различия.

Приоритет (и по чтению, и по записи) разбиения цикла с параметром j_ζ определяет возможность уменьшения r_ζ^9 без увеличения оценки коммуникационных издержек. Уменьшать r_ζ^9 требуется, если блоки вычислений являются слишком большими для эффективной реализации на мультипроцессорах. Целая величина приоритета предоставляет такую возможность, а нецелая – не предоставляет. Нецелый приоритет (вида $\tau + 0,5$, где $\tau \in Z, \tau \geq 0$) отражает наличие множителя Q_ζ^9 в оценке объема коммуникаций. Уменьшение r_ζ^9 увеличит Q_ζ^9 (в силу $M = O(r_\zeta^9 Q_\zeta^9)$) и, следовательно, оценку числа коммуникационных операций. Реальное число коммуникаций также может увеличиться. Например, нарушение условия $|\eta_\zeta^{\alpha,\beta}| < r_\zeta^9$ по причине малости r_ζ^9 во втором случае теоремы приведет к коммуникациям на каждой итерации (вместо нескольких единиц) из r_ζ^9 итераций j_ζ (условие (5) леммы 1) в каждой из Q_ζ^9 частей. В случае целого приоритета Q_ζ^9 не входит в оценку объема коммуникационных операций и, следовательно, есть возможность уменьшения r_ζ^9 без роста оценки. При ранжировании (установлении соотношений размеров относительно друг друга) параметров размера блоков вычислений параллельного алгоритма, реализуемого на графическом процессоре, следует учитывать в оценках объема коммуникационных операций только слагаемые с множителем Q_ζ^9 .

Пример. Рассмотрим алгоритм прямого хода метода Гаусса решения систем линейных алгебраических уравнений:

```

do  $k = 1, n - 1$ 
  do  $i = k + 1, n$ 
    do  $j = k + 1, n + 1$ 
       $a(i, j) = a(i, j) - \frac{a(i, k)}{a(k, k)} a(k, j)$ 
    enddo
  enddo
enddo

```

Матрицы $\Phi_{\alpha, \beta}$ и векторы $\varphi^{\alpha, \beta}$ в функциях зависимостей будем для наглядности помечать элементами массивов, фигурирующими на порождающих зависимости вхождениях. Например, матрицу $\Phi_{\alpha, \beta}$ и вектор $\varphi^{\alpha, \beta}$ для вхождения $a(i, k)$ обозначим $\Phi_{a(i, j), a(i, k)}$ и $\varphi^{a(i, j), a(i, k)}$.

Для второго вхождения массива a в оператор $a(i, j)$ – использование прежнего значения обновляемого элемента) имеем

$$F_{1,1,2} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \Phi_{a(i, j), a(i, j)} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \varphi^{a(i, j), a(i, j)} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \rho_{1,1,2}^{\Phi} = \text{rank} \begin{pmatrix} F_{1,1,2} \\ \Phi_{a(i, j), a(i, j)} \end{pmatrix} = 3.$$

Объемы коммуникационных операций $\omega_{1,1,2}^{\zeta, R}$ и $\omega_{1,1,2}^{\zeta, W}$, $\zeta = 1, 2, 3$, оцениваются величинами $\omega_{1,1,2}^{1, R} = \omega_{1,1,2}^{1, W} = O(Q_1 M^{\rho_{1,1,2}^{\Phi}}) = O(Q_1 M^2)$ (второй случай теоремы), $\omega_{1,1,2}^{2, R} = \omega_{1,1,2}^{2, W} = 0$ (первый случай), $\omega_{1,1,2}^{3, R} = \omega_{1,1,2}^{3, W} = 0$ (первый случай).

Для вхождения $a(i, k)$ (использование столбца, содержащего ведущий элемент) получим

$$F_{1,1,3} = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \Phi_{a(i, j), a(i, k)} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \varphi^{a(i, j), a(i, k)} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \rho_{1,1,3}^{\Phi} = \text{rank} \begin{pmatrix} F_{1,1,3} \\ \Phi_{a(i, j), a(i, k)} \end{pmatrix} = 2,$$

$\omega_{1,1,3}^{1, R} = \omega_{1,1,3}^{1, W} = O(Q_1 M^{\rho_{1,1,3}^{\Phi}}) = O(Q_1 M)$ (второй случай), $\omega_{1,1,3}^{2, R} = \omega_{1,1,3}^{2, W} = 0$ (первый случай), $\omega_{1,1,3}^{3, R} = O(Q_3 M^{\rho_{1,1,3}^{\Phi}}) = O(Q_3 M^2)$ и $\omega_{1,1,3}^{3, W} = O(M^{\rho_{1,1,3}^{\Phi}}) = O(M^2)$ (четвертый случай).

Для вхождения $a(k, j)$ (использование строки, содержащей ведущий элемент)

$$F_{1,1,4} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \Phi_{a(i, j), a(k, j)} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \varphi^{a(i, j), a(k, j)} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \rho_{1,1,4}^{\Phi} = \text{rank} \begin{pmatrix} F_{1,1,4} \\ \Phi_{a(i, j), a(k, j)} \end{pmatrix} = 2,$$

$\omega_{1,1,4}^{1, R} = \omega_{1,1,4}^{1, W} = O(Q_1 M^{\rho_{1,1,4}^{\Phi}}) = O(Q_1 M)$ (второй случай), $\omega_{1,1,4}^{2, R} = O(Q_2 M^{\rho_{1,1,4}^{\Phi}}) = O(Q_2 M^2)$ и $\omega_{1,1,4}^{2, W} = O(M^{\rho_{1,1,4}^{\Phi}}) = O(M^2)$ (четвертый случай), $\omega_{1,1,4}^{3, R} = \omega_{1,1,4}^{3, W} = 0$ (первый случай).

Для вхождения $a(k, k)$ (использование ведущего элемента)

$$F_{1,1,5} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \Phi_{a(i, j), a(k, k)} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \varphi^{a(i, j), a(k, k)} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \rho_{1,1,5}^{\Phi} = \text{rank} \begin{pmatrix} F_{1,1,5} \\ \Phi_{a(i, j), a(k, k)} \end{pmatrix} = 1,$$

$\omega_{1,1,5}^{1, R} = \omega_{1,1,5}^{1, W} = O(Q_1 M^{\rho_{1,1,5}^{\Phi}}) = O(Q_1)$ (второй случай), $\omega_{1,1,5}^{2, R} = O(Q_2 M^{\rho_{1,1,5}^{\Phi}}) = O(Q_2 M)$, $\omega_{1,1,5}^{2, W} = O(M^{\rho_{1,1,5}^{\Phi}}) = O(M)$, $\omega_{1,1,5}^{3, R} = O(Q_3 M^{\rho_{1,1,5}^{\Phi}}) = O(Q_3 M)$ и $\omega_{1,1,5}^{3, W} = O(M^{\rho_{1,1,5}^{\Phi}}) = O(M)$ (четвертый случай).

Объемы операций чтения, операций записи и суммарный объем коммуникационных операций для первой, второй и третьей координат оцениваются величинами $\omega_1^R = O(Q_1 M^2)$, $\omega_1^W = O(Q_1 M^2)$, $\omega_2^R = O(Q_2 M^2)$, $\omega_2^W = O(M^2)$, $\omega_3^R = O(Q_3 M^2)$, $\omega_3^W = O(M^2)$, $\omega_1 = O(Q_1 M^2)$, $\omega_2 = O(Q_2 M^2)$, $\omega_3 = O(Q_3 M^2)$. Первая, вторая и третья координаты имеют приоритеты для разбиения по чтению $z_1^R = 2, 5$, $z_2^R = 2, 5$, $z_3^R = 2, 5$, по записи $z_1^W = 2, 5$, $z_2^W = 2$, $z_3^W = 2$, по чтению и записи $z_1 = 2, 5$, $z_2 = 2, 5$, $z_3 = 2, 5$.

Для ранжирования параметров r_1 , r_2 и r_3 размера блоков вычислений графического процессора найдем объемы (для различия с вычисленными выше объемами будем помечать их индексом rank) операций чтения, операций записи и суммарный объем коммуникационных операций для первой, второй и третьей координат, учитывая только слагаемые с множителем Q_i : $\omega_1^{R,\text{rank}} = O(Q_1 M^2)$, $\omega_1^{W,\text{rank}} = O(Q_1 M^2)$, $\omega_2^{R,\text{rank}} = O(Q_2 M^2)$, $\omega_2^{W,\text{rank}} = 0$, $\omega_3^{R,\text{rank}} = O(Q_3 M^2)$, $\omega_3^{W,\text{rank}} = 0$, $\omega_1^{\text{rank}} = O(Q_1 M^2)$, $\omega_2^{\text{rank}} = O(Q_2 M^2)$, $\omega_3^{\text{rank}} = O(Q_3 M^2)$. Первая, вторая и третья координаты имеют приоритеты для разбиения по чтению $z_1^{R,\text{rank}} = 2,5$, $z_2^{R,\text{rank}} = 2,5$, $z_3^{R,\text{rank}} = 2,5$, по записи $z_1^{W,\text{rank}} = 2,5$, $z_2^{W,\text{rank}} = -1$, $z_3^{W,\text{rank}} = -1$, по чтению и записи $z_1^{\text{rank}} = 2,5$, $z_2^{\text{rank}} = 2,5$, $z_3^{\text{rank}} = 2,5$. При этом отсутствует возможность уменьшения r_1 без роста оценки коммуникационных издержек и по чтению, и по записи (нецелый приоритет $z_1^{R,\text{rank}} = z_1^{W,\text{rank}} = 2,5$); имеется возможность уменьшения r_2 и r_3 без роста оценки коммуникационных операций записи (целый приоритет $z_2^{W,\text{rank}} = z_3^{W,\text{rank}} = -1$). Поэтому уменьшать блоки вычислений следует в первую очередь за счет параметров размера r_2 и r_3 .

Таким образом, в работе сформулированы и доказаны утверждения, позволяющие оценить объем коммуникационных операций, порождаемых разбиением множества итераций. Определены приоритеты параметров циклов, позволяющие ранжировать параметры размера тайлов первого уровня для минимизации объема коммуникационных операций. Приведено обоснование старшинства приоритетов на основании числа обращений к глобальной памяти графического процессора.

Список использованной литературы

1. *Baskaran, M.* Automatic C-to-CUDA code generation for affine programs / M. Baskaran, J. Ramanujam, P. Sadayappan // Proceedings of the Compiler Construction, 19th International Conference. Part of the Joint European Conferences on Theory and Practice of Software. – Paphos, Cyprus, March 2010.
2. *Xue, J.* Time-minimal tiling when rise is larger than zero / J. Xue, W. Cai // Parallel Computing. – 2002. – Vol. 28, N 5. – P. 915–939.
3. *Kim, D. G.* Parameterized tiling for imperfectly nested loops / D. G. Kim, S. Rajopadhye // Technical Report CS-09-101, Colorado State University, Department of Computer Science, February 2009. – 21 p.
4. Automatic parallelization of tiled loop nests with enhanced fine-grained parallelism on GPUs / P. Di [et al.] // 41st International Conference on Parallel Processing. – Pittsburgh, PA, USA, September 2012. IEEE Computer Society, 2012. – P. 350–359.
5. *Bandishti, V.* Tiling stencil computations to maximize parallelism / V. Bandishti, I. Pananilath, U. Bondhugula // Proceedings of Supercomputing. – Los Alamitos, CA, USA. IEEE Computer Society Press, 2012. – P. 40:1–40:11.
6. *Воеводин, В. В.* Параллельные вычисления / В. В. Воеводин, Вл. В. Воеводин. – СПб.: БХВ-Петербург, 2002. – 608 с.
7. *Feautrier, P.* Some efficient solutions to the affine scheduling problem. Part 1 / P. Feautrier // International J. of Parallel Programming. – 1992. – Vol. 21, N 5. – P. 313–348.
8. Automatic transformations for communication-minimized parallelization and locality optimization in the polyhedral model / U. Bondhugula [et al.] // Lecture notes in computer science. – 2008. – N 4959. – P. 132–146.
9. *Лиходед, Н. А.* Оценка объема коммуникационных операций параллельного зернистого алгоритма / Н. А. Лиходед, М. А. Полещук // Междунар. конгресс по информатике: информационные системы и технологии CSIST'2013, 4–7 ноября 2013 г., Минск, Беларусь. Бел. гос. ун-т. – Минск, 2013. – С. 377–381.
10. *Лиходед, Н. А.* Характеристика локальности параллельных реализаций многомерных циклов / Н. А. Лиходед // Докл. НАН Беларуси. – 2010. – Т. 54, № 1. – С. 26–32.
11. *Адуцкевич, Е. В.* К распараллеливанию последовательных программ: распределение массивов между процессорами и структуризация коммуникаций / Е. В. Адуцкевич, Н. А. Лиходед, А. О. Сикорский // Кибернетика и системный анализ. – 2012. – Т. 48, № 1. – С. 144–163.

Поступило в редакцию 30.03.2015