${f B}$ ычислительная математика

Computational mathematics

УДК 519.67

УСЛОВИЯ ПРИВАТИЗАЦИИ ЭЛЕМЕНТОВ МАССИВА ПОТОКАМИ ВЫЧИСЛЕНИЙ

 $H. A. ЛИХОДЕД^{1}, M. A. ПОЛЕЩУ<math>K^{1}$

1)Белорусский государственный университет, пр. Независимости, 4, 220030, г. Минск, Беларусь

Множество операций параллельного алгоритма для реализации на графическом процессоре должно быть разбито на потоки (нити) вычислений. Потоки следует сгруппировать в блоки вычислений, выполняющиеся атомарно на потоковых процессорах, называемых также мультипроцессорами. Для хорошей производительности графического процессора важно, чтобы как можно больше данных умещались в быстрых регистровой и разделяемой памяти, иначе используются медленные глобальная и локальная память. Степень использования памяти с быстрым доступом отражает вычислительное свойство алгоритма, называемое локальностью. При реализации алгоритмов на многопроцессорных вычислительных устройствах применение локальности играет важнейшую роль для достижения высокой производительности. В данной работе сформулированы и доказаны необходимые условия и достаточные условия, использование которых позволяет получать потоки с приватизированными данными, т. е. такие потоки вычислений, что элемент массива используется только одним потоком, и поэтому его целесообразно разместить в регистре.

Ключевые слова: параллельные вычисления; графический процессор; тайлинг; приватизация элементов массива; регистры.

Благодарность. Работа выполнена в рамках государственной программы научных исследований Республики Беларусь «Конвергенция-2020» (подпрограмма «Методы математического моделирования сложных систем»).

Образец цитирования:

Лиходед НА, Полещук МА. Условия приватизации элементов массива потоками вычислений. Журнал Белорусского государственного университета. Математика. Информатика. 2018;3:59—67.

For citation:

Likhoded NA, Paliashchuk MA. Conditions for privatizing the elements of arrays by computing threads. *Journal of the Belarusian State University. Mathematics and Informatics.* 2018;3: 59–67. Russian

Авторы:

Николай Александрович Лиходед – доктор физико-математических наук, профессор; профессор кафедры вычислительной математики факультета прикладной математики и информатики.

Максим Александрович Полещук – ассистент кафедры вычислительной математики факультета прикладной математики и информатики.

Authors:

Nikolai A. Likhoded, doctor of science (physics and mathematics), full professor; professor at the department of computational mathematics, faculty of applied mathematics and computer science.

likhoded@bsu.bv

Maksim A. Paliashchuk, assistant at the department of computational mathematics, faculty of applied mathematics and computer science.

poleschuma@bsu.by

CONDITIONS FOR PRIVATIZING THE ELEMENTS OF ARRAYS BY COMPUTING THREADS

N. A. LIKHODED^a, M. A. PALIASHCHUK^a

^aBelarusian State University, 4 Niezaliežnasci Avenue, Minsk 220030, Belarus Corresponding author: N. A. Likhoded (likhoded@bsu.by)

The set of operations of the parallel algorithm for implementation on the GPU must be split into computation threads. The threads must be grouped into computation units that run atomically on stream processors, also called multiprocessors. For good GPU performance, it is important that as much data as possible can fit into fast register and shared memory, otherwise slow global and local memory are used. The degree of memory usage with fast access reflects the computational property of the algorithm, called locality. When implementing algorithms on multiprocessor computing devices, the use of locality plays a crucial role in achieving high performance. In this paper, necessary conditions and sufficient conditions have been formulated and proved, the use of which allows receiving threads with privatized data, i. e. it allows to receive such computation threads that the array element is used only by one thread and therefore it is advisable to place it in the register.

Key words: parallel computations; GPU; tiling; array privatization; registers.

Acknowledgements. The prepared report was sponsored by the government program of scientific research of the Republic of Belarus «Convergence-2020» (subprogram «Methods of mathematical modeling of complex systems»).

Введение

В качестве целевого компьютера для реализации алгоритмов будем рассматривать графические процессоры (GPU) — довольно мощные и в настоящее время широко используемые многоядерные устройства. При вычислениях на GPU быстрым является процесс обращения к регистрам (самый быстрый вид памяти), разделяемой памяти мультипроцессора и кешам, но не обращение к глобальной и локальной памяти GPU. Регистры мультипроцессора делятся поровну между потоками блока вычислений. Особенно эффективно используется регистровая память, если элемент массива приватизирован потоком. Под приватизацией понимается использование элемента массива в вычислениях только одного потока блока вычислений. Если элемент массива приватизирован, то он может быть размещен в регистрах, выделенных потоку. Цель данной работы — формулировка и доказательство необходимых условий и достаточных условий приватизации элементов массива потоками вычислений.

Степень использования памяти с быстрым доступом отражает вычислительное свойство алгоритма, называемое локальностью. При реализации алгоритмов на многопроцессорных устройствах локальность играет важнейшую роль для достижения высокой производительности [1; 2].

Отметим некоторые исследования локальности алгоритмов для реализации на GPU. В работе [2] на примере решения задачи поиска кратчайших путей показана особая для GPU эффективность использования свойства локальности. В работе [3] сформулированы и доказаны утверждения, позволяющие ранжировать параметры размера блоков на основе асимптотических оценок объема коммуникационных операций; в [4] подобные оценки получены с использованием более сложного (но и лучше приспособленного для автоматизации) математического аппарата. Известен способ получения более точной оценки количества элементов, к которым осуществляется доступ при выполнении операций блока вычислений, но практическое применение этого результата довольно трудоемкое и требует привлечения специализированных средств автоматизации [5; 6]. В работе [7] сформулированы необходимые условия использования элемента массива только одним потоком вычислений и, соответственно, размещения в регистрах.

Разбиение множества операций алгоритма на блоки и потоки вычислений

Пусть алгоритм задан гнездом вложенных циклов, в котором имеется Θ наборов выполняемых операторов. Под набором операторов будем понимать один или несколько выполняемых операторов, окруженных одним и тем же множеством циклов. Выполняемые операторы и наборы операторов линейно упорядочены расположением их в записи алгоритма. Обозначим: V^{θ} – область изменения параметров циклов, окружающих θ -й набор операторов, $1 \le \theta \le \Theta$; n^{θ} – размерность этой области, число циклов,

окружающих θ -й набор операторов; v_l – размерности массивов a_l . Также обозначим размеры блоков вычислений натуральными числами $r_1^{\theta}, \ldots, r_{n^{\theta}}^{\theta}$; r_{ζ}^{θ} – число значений параметра j_{ζ} , приходящихся на один блок θ -го набора операторов; Q_{ζ}^{θ} – число блоков θ -го набора операторов по координате с номером ζ (обозначение этой координаты — j_{ζ}^{gl}). Нумеровать блоки вычислений будем по каждой координате в пределах от 0 до Q_{ζ}^{θ} – 1, $1 \leq \zeta \leq n^{\theta}$. Блоки обозначим $V_{j^{gl}}^{gl}$, где $J^{gl} = \left(j_1^{gl}, \ldots, j_{n^{\theta}}^{gl}\right)$, $0 \leq j_{\zeta}^{gl} \leq Q_{\zeta}^{\theta} - 1$, $1 \leq \zeta \leq n^{\theta}$. Далее для простоты записи будем часто использовать обозначения без индекса θ .

Зададим в блоках вычислений потоки вычислений посредством выделения блоков второго уровня. Зададим размеры потоков натуральными числами $r_{1,2}, r_{2,2}, ..., r_{n^0,2}; r_{\zeta,2}$ – число значений параметра j_{ζ} ,

приходящихся на один поток;
$$Q_{1,2}, Q_{2,2}, ..., Q_{n^0,2}$$
 – число таких потоков; $Q_{\zeta,2} = \left[\frac{r_{\zeta}}{r_{\zeta,2}}\right]$.

В общем случае разбить множество операций алгоритма на блоки и потоки можно путем тайлинга – преобразования алгоритма для получения макроопераций [8].

Вхождением (a_l, S_{β}, q) будем называть q-е вхождение массива a_l в оператор S_{β} . Индексы элементов l-го массива, связанных с вхождением (a_l, S_{β}, q) , выражаются функцией вида

$$\overline{F}_{a_{l}, S_{\beta}, q}(J) = F_{a_{l}, S_{\beta}, q}J + f^{a_{l}, S_{\beta}, q}, J(j_{1}, ..., j_{n^{\theta}}) \in V^{\theta}, F_{a_{l}, S_{\beta}, q} \in \mathbb{Z}^{v_{l} \times n^{\theta}}, f^{a_{l}, S_{\beta}, q} \in \mathbb{Z}^{v_{l}}$$

Пример 1. Рассмотрим основную часть алгоритма Флойда – Уоршелла поиска кратчайших путей между всеми парами вершин графа:

```
do k = 1, n

do i = 1, n

do j = 1, n

a(i, j) = \min(a(i, j), a(i, k) + a(k, j))

enddo

enddo

enddo
```

В гнезде циклов имеется один выполняемый оператор S_1 (один набор операторов), и используется один массив a размерности 2. Область изменения параметров циклов (область итераций) $V^1 = \{(k,i,j) \in \mathbb{Z}^3 \mid 1 \le k \le n, \ 1 \le i \le n, \ 1 \le j \le n\}$ для оператора S_1 имеет размерность 3. Для матриц $F_{a,S_1,q}$ на вхождениях (a,S_1,q) имеем

$$F_{a,S_1,1} = F_{a,S_2,2} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \ F_{a,S_1,3} = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \ F_{a,S_1,4} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Известны блочные алгоритмы с 3D-блоками размером $r \times r \times r$ для реализации на GPU [9; 10]. Заметим, что одинаковые размеры блока имеют существенное значение, а не выбраны для простоты. Пусть k^{gl} , i^{gl} , j^{gl} — номера частей, на которые при формировании блоков разбиваются области значений параметров k, i, j циклов. Блоки $V_{k^{gl}}$ i^{gl} i^{gl} i^{gl} i^{gl} имеют следующий вид:

do
$$k = 1 + k^{gl}r$$
, $\min((k^{gl} + 1)r, n)$
do $i = 1 + i^{gl}r$, $\min((i^{gl} + 1)r, n)$
do $j = 1 + j^{gl}r$, $\min((j^{gl} + 1)r, n)$
 $a(i, j) = \min(a(i, j), a(i, k) + a(k, j))$
enddo
enddo
enddo

Обозначим через $r_{k,\,2},\,r_{i,\,2}$ и $r_{j,\,2}$ размеры блоков второго уровня. Числа $Q_{k,\,2} = \left[\frac{r}{r_{k,\,2}}\right],\,Q_{i,\,2} = \left[\frac{r}{r_{i,\,2}}\right]$

и $Q_{j,2} = \left[\frac{r}{r_{j,2}}\right]$ задают количество итераций в новых циклах с параметрами $k^{g/2}$, $i^{g/2}$ и $j^{g/2}$ соответственно.

Пусть $r_{k,\,2}=r$, тогда $Q_{k,\,2}=1$. Блок $V_{k^{gl},\,i^{gl},\,j^{gl}}$ с выделенными потоками $\operatorname{Thr}\left(i^{gl2},\,j^{gl2}\right)$ имеет следующее представление:

do
$$i^{gl2} = 0$$
, $Q_{i,2} - 1$
do $j^{gl2} = 0$, $Q_{j,2} - 1$
Thr (i^{gl2}, j^{gl2}) :
do $k = 1 + k^{gl}r$, min $((k^{gl} + 1)r, n)$
do $i = 1 + i^{gl}r + i^{gl2}r_{i,2}$, min $(i^{gl}r + (i^{gl2} + 1)r_{i,2}, (i^{gl2} + 1)r, n)$
do $j = 1 + j^{gl}r + j^{gl2}r_{j,2}$, min $(j^{gl}r + (j^{gl2} + 1)r_{j,2}, (j^{gl} + 1)r, n)$
 $a(i, j) = \min(a(i, j), a(i, k) + a(k, j))$
enddo
enddo
enddo
enddo
enddo

Необходимое условие приватизации элементов массива

Сформулируем и докажем условия, налагаемые на матрицу $F_{a_l, S_{\beta}, q}$, без выполнения которых невозможна приватизация элементов массива, используемых на вхождении (a_l, S_{β}, q) , в любом из трех возможных способов задания потоков вычислений (потоки могут помечаться как узлы одномерной, двумерной или трехмерной сетки).

Введем в рассмотрение множество $Z_s = \{\zeta_1, ..., \zeta_m\} \subset \{1, ..., n^\theta\}$, составленное из одного (m=1), двух (m=2) или трех (m=3) произвольных элементов множества $\{1, ..., n^\theta\}$, и множество $Z_t = \{\zeta_{m+1}, ..., \zeta_{n^\theta}\} = \{1, ..., n^\theta\} \setminus Z_s$. Для удобства будем считать множества Z_s и Z_t упорядоченными: $\zeta_i < \zeta_j$, если i < j (отдельно для любых ζ_i , $\zeta_j \in Z_s$ и ζ_i , $\zeta_j \in Z_t$).

Полагаем, что $j_{\zeta_1}^{g/2}, \ldots, j_{\zeta_m}^{g/2}$ задают координаты потоков, и если $\zeta \in Z_s$, то $1 \le r_{\zeta, 2} < r_{\zeta}$, а если $\zeta \in Z_t$, то $r_{\zeta, 2} = r_{\zeta}$. Таким образом, m координат задают m-мерное пространство потоков, а остальные $n^{\theta} - m$ координат – итерации, выполняемые потоками в лексикографическом порядке.

Введем обозначения: T_s — матрица, строки которой составлены из векторов $e_{\zeta}^{(n^0)}$, $\zeta \in Z_s$, где $e_{\zeta}^{(n^0)}$ — векторо-строка размером n^0 , у которого координата с номером ζ равна 1, а остальные координаты нулевые;

$$\rho_{a_{l},\,S_{\beta},\,q} = \mathrm{rank}\,F_{a_{l},\,S_{\beta},\,q},\ \, \rho_{a_{l},\,S_{\beta},\,q}^{s} = \mathrm{rank}\binom{F_{a_{l},\,S_{\beta},\,q}}{T_{s}};\ u_{i} - \text{базисные векторы (размерности } n^{\theta})\ \text{подпространства}$$

 $\ker F_{a_i, S_{\beta}, q}$ пространства $Z^{n^{\theta}}$, $1 \le i \le n^{\theta} - \rho_{a_i, S_{\beta}, q}$. Если $\rho_{a_i, S_{\beta}, q} = n^{\theta}$, то подпространство $\ker F_{a_i, S_{\beta}, q}$ имеет нулевую размерность.

Теорема 1. Пусть зафиксирован блок вычислений, потоки вычислений ${
m Thr}ig(j_{\zeta_1}^{gl2},...,j_{\zeta_m}^{gl2}ig)$ помечены как узлы одномерной, или двумерной, или трехмерной сетки и для вхождения $ig(a_l,S_{eta},qig)$ массива a_l в оператор S_{eta} выполнено условие

$$\rho_{a_l, S_B, q} = \rho_{a_l, S_B, q}^s. \tag{2}$$

Тогда каждый элемент массива a_l используется на вхождении $\left(a_l,\,S_{\beta},\,q\right)$ только одним потоком. Если $\rho_{a_l,\,S_{\beta},\,q} < \rho^s_{a_l,\,S_{\beta},\,q}$, то элемент массива используется, вообще говоря, более чем одним потоком.

Доказательство. Сперва рассмотрим случай $r_{\zeta,\,2}=1,\,\zeta\in Z_s$. Тогда один поток $\operatorname{Thr}\left(j_{\zeta_1},\,...,\,j_{\zeta_m}\right)^{\mathrm{T}}=$ составляют все вычисления блока при фиксированных значениях $j_{\zeta_1},\,...,\,j_{\zeta_m}$, причем $\left(j_{\zeta_1},\,...,\,j_{\zeta_m}\right)^{\mathrm{T}}=T_s\left(j_1,\,...,\,j_{\eta^0}\right)^{\mathrm{T}}.$

Одним из результатов работы [11] (частный случай следствия 2) является следующее утверждение: для каждого элемента массива, связанного с вхождением (a_i, S_β, q) , значения координат потоков, в которых используется этот элемент, отличаются линейными комбинациями векторов $T_s u_i$, $1 \le i \le n^\theta - \rho_{a_i, S_\beta, q}$.

Если $\rho_{a_l, S_{\beta}, q} = \rho_{a_l, S_{\beta}, q}^s$, то $\ker \begin{pmatrix} F_{a_l, S_{\beta}, q} \\ T_s \end{pmatrix} = \ker F_{a_l, S_{\beta}, q}$, поэтому $T_s u_i = 0$ для всех u_i . Следовательно, в этом слу-

чае каждый элемент массива используется только одним потоком. Если $\rho_{a_l, S_{\beta}, q} < \rho_{a_l, S_{\beta}, q}^s$, то $\ker \begin{pmatrix} F_{a_l, S_{\beta}, q} \\ T_s \end{pmatrix} \neq$

 $\neq \ker F_{a_i, S_\beta, q}$, поэтому для каких-то u_i имеет место $T_s u_i \neq 0$. Тогда элемент массива используется (если только этот элемент не используется однократно на границе блока вычислений) более чем одним потоком.

Пусть теперь $1 < r_{\zeta,2} < r_{\zeta}$ для каких-то (возможно, всех) $\zeta \in Z_s$. Тогда один поток $\mathrm{Thr}\left(j_{\zeta_1}^{gl2},\ldots,j_{\zeta_m}^{gl2}\right)$ составляют вычисления с более чем одним значением таких j_{ζ} . Если элемент массива используется только одним потоком при $r_{\zeta,2}=1$, то тем более он будет использоваться только одним потоком при $r_{\zeta,2}>1$. Если элемент массива используется более чем одним потоком при $r_{\zeta,2}=1$, то при $r_{\zeta,2}>1$ он будет использоваться меньшим числом потоков, но с учетом $r_{\zeta,2}< r_{\zeta}$ все же более чем одним (исключение, как и при $r_{\zeta,2}=1$, возможно на границе блока вычислений).

Частные случаи теоремы 1 (m = 1, m = 2) сформулированы в работе [7].

Пример 1 (продолжение). Пусть блоки имеют вид (1). Наиболее интересными с точки зрения эффективного использования быстрой регистровой памяти являются такие варианты выбора матрицы T_s , при которых условие (2) выполняется для вхождений $(a, S_1, 1)$ и $(a, S_1, 2)$: на этих вхождениях осуществляется считывание и запись данного a(i, j) (на других вхождениях – только считывание a(i, k) или a(k, j)). Нетрудно видеть, что условия $\rho_{a, S_1, 1} = \rho_{a, S_1, 1}^s$, $\rho_{a, S_1, 2} = \rho_{a, S_1, 2}^s$ выполняются, если матрица T_s не содержит строки $e_1^{(3)} = (1 \ 0 \ 0)$. Например, если $Z_s = \{2, 3\}$, то

$$\rho_{a, S_1, 1} = \operatorname{rank} F_{a, S_1, 1} = \operatorname{rank} \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = 2, \ \rho_{a, S_1, 1}^s = \operatorname{rank} \begin{pmatrix} F_{a, S_1, 1} \\ T_s \end{pmatrix} = \operatorname{rank} \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = 2.$$

Пример 2. Основную часть алгоритма численного решения двумерной задачи Дирихле для уравнения Пуассона методом последовательной верхней релаксации можно представить (после аффинного преобразования итерационной области, применяемого для возможности получать блочные версии) в следующем виде:

```
do t = 1, T

do i = t + 1, t + N_x - 1

do j = t + 1, t + N_y - 1

y(i - t, j - t) = F(y(i - t - 1, j - t), y(i - t, j - t - 1), y(i - t, j - t), y(i - t + 1, j - t), y(i - t, j - t + 1))

enddo

enddo

enddo
```

Здесь T — некоторое фиксированное число итераций метода релаксации, N_x и N_y характеризуют число частей, на которые разбивается итерационная область.

Для любого вхождения (y, S_1, q) и для любого выбора матрицы T_s имеем

$$F_{a_{l}, S_{\beta}, q} = \begin{pmatrix} -1 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix}, \ \rho_{a_{l}, S_{\beta}, q} = 2, \ \rho_{a_{l}, S_{\beta}, q}^{s} = 3$$

(ранг матрицы $F_{a_l, S_\beta, q}$ равен 2, но после добавления в матрицу хотя бы одной из строк (1 0 0), (0 1 0), (0 0 1) ранг полученной матрицы станет равен 3). Поэтому для любого варианта организации потоков выполняется $\rho_{a_l, S_\beta, q} < \rho_{a_l, S_\beta, q}^s$, каждый элемент массива используется более чем одним потоком, и элементы массива y не могут размещаться в регистрах.

Достаточные условия приватизации элементов массива

Условие (2) гарантирует использование элемента массива a_l на вхождении (a_l, S_β, q) только одним потоком. Для исследования возможности приватизации элемента массива надо принять во внимание и другие вхождения элементов этого массива в операторы данного блока вычислений.

Обозначим через $\overline{F}_{a_l,\,S_{eta},\,q}ig(V_{J^{gl}}^{m{ heta}}ig)$ образ множества $V_{J^{gl}}^{m{ heta}}$ после применения отображения $\overline{F}_{a_l,\,S_{m{eta},\,q}}$. По определению

$$\overline{F}_{a_{l}, S_{\beta}, q}\left(V_{J^{\mathcal{B}^{l}}}^{\theta}\right) = \left\{ F \in \mathbb{Z}^{\mathsf{v}_{l}} \middle| F = \overline{F}_{a_{l}, S_{\beta}, q}(J), J \in V_{J^{\mathcal{B}^{l}}}^{\theta} \right\}.$$

Теорема 2. Пусть зафиксирован блок вычислений $V_{J^{gl}}^{\theta}$, потоки вычислений $\mathrm{Thr} \left(j_{\zeta_1}^{gl2}, \ldots, j_{\zeta_m}^{gl2} \right)$ помечены как узлы одномерной, или двумерной, или трехмерной сетки и для вхождения $\left(a_l, S_{\beta}, q \right)$ массива a_l в оператор S_{β} выполнено условие (2). Для того чтобы каждый элемент массива a_l , возникающий на вхождении $\left(a_l, S_{\beta}, q \right)$, использовался только одним потоком, достаточно выполнения условий

$$\bar{F}_{a_l, S_{\beta}, q}(V_{J^{gl}}^{\theta}) \cap \bar{F}_{a_l, S_{\gamma}, y}(V_{J^{gl}}^{\theta}) = \varnothing$$

или

$$F_{a_l, S_{\gamma}, y} = F_{a_l, S_{\beta}, q}, \ f^{a_l, S_{\gamma}, y} = f^{a_l, S_{\beta}, q},$$

где (a_l, S_{γ}, y) есть любое, отличное от (a_l, S_{β}, q) вхождение массива a_l в операторы блока вычислений V_{rsl}^{θ} .

лений $V_{J^{s'}}^{\theta}$. Доказательство. Пусть имеет место условие (2) и пусть сначала $r_{\zeta,\,2}=1,\,\zeta\in Z_s$. Тогда один поток ${
m Thr}ig(j_{\zeta_1},\,...,\,j_{\zeta_m}ig)$ составляют все вычисления блока при фиксированных значениях $j_{\zeta_1},\,...,\,j_{\zeta_m}$ и на вхождении $a_l,\,S_{\beta},\,a_l$ каждый элемент массива a_l используется только одним потоком.

Если $\overline{F}_{a_l, S_{\beta}, q}(V_{J^{gl}}^{\theta}) \cap \overline{F}_{a_l, S_{\gamma}, y}(V_{J^{gl}}^{\theta}) = \emptyset$, то на вхождении (a_l, S_{β}, q) и на любом другом вхождении (a_l, S_{γ}, y) используются разные элементы массива, поэтому другие вхождения не влияют на использование массива a_l на вхождении (a_l, S_{β}, q) только одним потоком.

Пусть $F_{a_l, S_\gamma, y} = F_{a_l, S_\beta, q}$, $f^{a_l, S_\gamma, y} = f^{a_l, S_\beta, q}$. Покажем, что на вхождениях (a_l, S_β, q) и (a_l, S_γ, y) один и тот же элемент массива может появиться только в одном потоке, т. е. если $\overline{F}_{a_l, S_\beta, q}(J) = \overline{F}_{a_l, S_\gamma, y}(I)$, $I, J \in V_{ls}^{\theta}$, то вычисления итераций I и J выполняются одним потоком.

Так как $F_{a_i, S_{\beta}, q}J + f^{a_i, S_{\beta}, q} = F_{a_i, S_{\gamma}, y}I + f^{a_i, S_{\gamma}, y}$, $F_{a_i, S_{\beta}, q} \left(J - I\right) = 0$, то $J - I \in \ker F_{a_i, S_{\beta}, q}$. Поэтому векторы I и J отличаются линейными комбинациями базисных векторов u_i , $1 \le i \le n^\theta - \rho_{a_i, S_{\beta}, q}$, подпространства $\ker F_{a_i, S_{\beta}, q}$. Тогда координаты потоков $T_s I$ и $T_s J$, которые выполняют операции итераций I и J, отличаются линейными комбинациями векторов $T_s u_i$. Так как $\rho_{a_i, S_{\beta}, q} = \rho^s_{a_i, S_{\beta}, q}$, то $T_s u_i = 0$ для всех u_i , вычисления итераций I и J выполняются одним потоком.

Пусть теперь $1 < r_{\zeta,\,2} < r_\zeta$ для каких-то (или всех) $\zeta \in Z_s$. Тогда один поток $\operatorname{Thr} \left(j_{\zeta_1}^{g/2}, \ldots, j_{\zeta_m}^{g/2}\right)$ составляют вычисления с более чем одним значением таких j_ζ . Если элемент массива используется только одним потоком при $r_{\zeta,\,2} = 1$, то тем более и при $r_{\zeta,\,2} > 1$ он будет использоваться лишь одним потоком. **Пример 1** (окончание). Рассмотрим блоки вычислений $V_{k^{gl},\,i^{gl},\,j^{gl}}$, для которых $i^{gl} \neq k^{gl}, j^{gl} \neq k^{gl}$. Покажем,

что если
$$T_s = \begin{pmatrix} 0 & 1 & 0 \end{pmatrix}$$
, или $T_s = \begin{pmatrix} 0 & 0 & 1 \end{pmatrix}$, или $T_s = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$, то на вхождениях $\begin{pmatrix} a, S_1, 1 \end{pmatrix}$ и $\begin{pmatrix} a, S_1, 2 \end{pmatrix}$

происходит приватизация потоками данных a(i, j). В предыдущем разделе установлено, что при таком выборе матрицы T_s условие (2) выполнено. Осталось показать, что выполнены и остальные условия теоремы 2. Имеем

$$\begin{split} \overline{F}_{a, S_{1}, 1}(k, i, j) &= \overline{F}_{a, S_{1}, 2}(k, i, j) = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} (k, i, j)^{\mathsf{T}} = (i, j)^{\mathsf{T}}, \\ \overline{F}_{a, S_{1}, 1}(V_{k^{gl}, i^{gl}, j^{gl}}) &= \overline{F}_{a, S_{1}, 2}(V_{k^{gl}, i^{gl}, j^{gl}}) = \\ &= \left\{ (F_{1}, F_{2}) \in \mathbb{Z}^{2} \middle| (F_{1}, F_{2}) = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} (k, i, j)^{\mathsf{T}} = (i, j), (k, i, j) \in V_{k^{gl}, i^{gl}, j^{gl}} \right\} = \\ &= \left\{ (F_{1}, F_{2}) \in \mathbb{Z}^{2} \middle| 1 + i^{gl}r \leq F_{1} \leq (i^{gl} + 1), 1 + j^{gl}r \leq F_{2} \leq (j^{gl} + 1) \right\}. \end{split}$$

Для вхождений $(a, S_1, 3), (a, S_1, 4),$ на которых осуществляется считывание данных a(i, k), a(k, j),имеем соответственно:

$$\begin{split} \overline{F}_{a,\,S_1,\,3}\big(k,\,i,\,j\big) &= \big(i,\,k\big)^{\mathsf{T}},\\ \overline{F}_{a,\,S_1,\,3}\Big(V_{k^{gl},\,i^{gl},\,j^{gl}}\Big) &= \Big\{\big(F_1,\,F_2\big) \in \mathbb{Z}^2 \,\Big|\, 1 + i^{gl}r \leq F_1 \leq \big(i^{gl}+1\big),\, 1 + k^{gl}r \leq F_2 \leq \big(k^{gl}+1\big)\Big\};\\ \overline{F}_{a,\,S_1,\,4}\big(k,\,i,\,j\big) &= \big(k,\,j\big)^{\mathsf{T}},\\ \overline{F}_{a,\,S_1,\,4}\Big(V_{k^{gl},\,i^{gl},\,j^{gl}}\Big) &= \Big\{\big(F_1,\,F_2\big) \in \mathbb{Z}^2 \,\Big|\, 1 + k^{gl}r \leq F_1 \leq \big(k^{gl}+1\big),\, 1 + j^{gl}r \leq F_2 \leq \big(j^{gl}+1\big)\Big\}. \end{split}$$

Так как $i^{gl} \neq k^{gl}$, $j^{gl} \neq k^{gl}$, то $\overline{F}_{a,\,S_i,\,1}\Big(V_{k^{gl}\ i^{gl}\ i^{gl}}\Big) \cap \overline{F}_{a,\,S_i,\,3}\Big(V_{k^{gl}\ i^{gl}\ j^{gl}}\Big) = \varnothing$ (пересечений нет по второй координахи нате), $\bar{F}_{a,\,S_1,\,1}\!\!\left(V_{k^{g_l},\,i^{g_l},\,j^{g_l}}\right)\cap \bar{F}_{a,\,S_1,\,4}\!\left(V_{k^{g_l},\,i^{g_l},\,j^{g_l}}\right) = \varnothing$ (пересечений нет по первой координате). Следовательно, выполняются достаточные условия приватизации.

Замечание 1. В рассмотренном примере множества $\bar{F}_{a,\,S_{i},\,q}\!\left(V_{k^{g_{i}},\,i^{g_{i}},\,j^{g_{i}}}\right)$ получить довольно просто. Для построения $ar{F}_{a_l, S_8, q}ig(V_{J^{sl}}^{ heta}ig)$ в общем случае можно воспользоваться методом получения отображений выпуклых ограниченных многогранников [12; 13].

3амечание 2. Пусть, как и ранее, $\operatorname{Thr}\left(j_{\zeta_1}^{gl2},\ldots,j_{\zeta_m}^{gl2}\right)$ – потоки вычислений, для вхождения $\left(a_l,\,S_{\mathsf{B}},\,q\right)$ массива a_l в оператор S_{β} выполнено условие (2), а для любого другого вхождения (a_l, S_{γ}, y) справедливо $F_{a_l, S_q, y} = F_{a_l, S_\beta, q}$, m строк матрицы $F_{a_l, S_\beta, q}$ совпадают со строками матрицы T_s и соответствующие координаты вектора $f^{a_l,\,S_\beta,\,q}-f^{a_l,\,S_\gamma,\,y}$ равны нулю. Тогда каждый элемент массива a_l , возникающий на вхождениях $(a_l, S_{\beta}, q), (a_l, S_{\gamma}, y)$, используется только одним потоком.

Действительно, покажем, что если $\overline{F}_{a_l,\,S_{\mathfrak{g}},\,q}\big(J\big) = \overline{F}_{a_l,\,S_{\mathfrak{q}},\,y}\big(I\big),\ I,\,J\in V_{J^{\mathfrak{g}l}}^{\,\theta},$ то $T_sI=T_sJ.$ Так как $F_{a_l,\,S_{\mathfrak{g}},\,q}J+T_sJ$ $+f^{a_{l},\,S_{\beta},\,q}=F_{a_{l},\,S_{\gamma},\,y}I+f^{a_{l},\,S_{\gamma},\,y},\quad F_{a_{l},\,S_{\gamma},\,y}=F_{a_{l},\,S_{\beta},\,q},\quad \text{то}\quad F_{a_{l},\,S_{\beta},\,q}\Big(J-I\Big)=f^{a_{l},\,S_{\beta},\,q}-f^{a_{l},\,S_{\gamma},\,y}.$ Отсюда, поскольку m строк матрицы $F_{a_l, S_{\beta}, q}$ совпадают со строками матрицы T_s и соответствующие координаты вектора $f^{a_l, S_{\beta}, q} - f^{a_l, S_{\gamma}, y}$ равны нулю, получим $T_s(J - I) = 0$.

Достаточные условия, приведенные в замечании 2, можно применить, например, в случае итераций (k, i, j), оператора вида $a(i, j) = a(i, j - 1) + ..., T_s = \begin{pmatrix} 0 & 1 & 0 \end{pmatrix}$ (потоки определяет координата i).

Заключение

В работе сформулированы и доказаны необходимые условия и достаточные условия приватизации элементов массива потоками вычислений при реализации алгоритма на графическом процессоре. Приватизация элементов массива позволяет разместить их в самой быстрой памяти процессора – регистрах.

Направления дальнейших исследований: формализация условий, при которых также происходит эффективное использование быстрой памяти графического процессора — бродкаст (потоки запрашивают одновременно один и тот же элемент массива) и пространственная локальность (потоки используют близко расположенные в памяти данные); применение сформулированных условий для автоматизации получения алгоритмов, реализуемых на графических процессорах; использование результатов работы при разработке параллельных алгоритмов для решения прикладных задач.

Библиографические ссылки

- 1. Воеводин ВлВ, Воеводин ВадВ. Спасительная локальность суперкомпьютеров. Открытые системы. 2013;9:12–15.
- 2. Buluc A, Gilberta JR, Budak C. Solving path problems on the GPU. Parallel Computing. 2010;36(5-6):241-253.
- 3. Лиходед НА, Полещук МА. Оценка локальности параллельных алгоритмов, реализуемых на графических процессорах. Вестник Южно-Уральского государственного университета. Серия: Вычислительная математика и информатика. 2016:5(3):96–111. DOI: 10.14529/cmse160307.
- 4. Лиходед НА, Полещук МА. Метод ранжирования параметров размера блоков вычислений параллельного алгоритма. Доклады Национальной академии наук Беларуси. 2015;59(4):25–33.
- 5. Kandemir M, Ramanujam J, Irwin M, Narayanan V, Kadayif I, Parikh A. A compiler based approach for dynamically managing scratch-pad memories in embedded systems. *IEEE Transactions on Computer-Aided Design*. 2004;23(2):243–260.
- 6. Baskaran M, Bondhugula U, Krishnamoorthy S, Ramanujam J, Rountev A, Sadayappan P. Automatic data movement and computation mapping for multi-level parallel architectures with explicitly managed memories. In: *Proceedings of the 13th ACM SIGPLAN Symposium on Principles and practice of parallel programming; 2008 February 20–23; Salt Lake City, USA.* New York: ACM; 2008. p. 1–10. DOI: 10.1145/1345206.1345210.
- 7. Полещук МА, Лиходед НА. Приватизация элементов массивов потоками вычислений. В: CSIST'2016. Международный конгресс по информащие: информационные системы и технологии; 24–27 октября 2016 г.; Минск, Беларусь. Минск: БГУ; 2016 г. 883–888
- 8. Baskaran M, Ramanujam J, Sadayappan P. Automatic C-to-CUDA code generation for affine programs. In: *Compiler Construction*. 19th International Conference. Part of the Joint European Conferences on Theory and Practice of Software; 2010 March 20–28; Paphos, Cyprus. Berlin, Heidelberg: Springer-Verlag; 2010. DOI: 10.1007/978-3-642-11970-5_14.
- 9. Katz GJ, Kider J. All-pairs shortest-paths for large graphs on the GPU. In: *Proceedings of the 23rd ACM SIGGRAPH/EURO-GRAPHICS symposium on Graphics hardware; 2008 June 20–21; Sarajevo, Bosnia and Herzegovina.* Geneve: Eurographics Association; 2008. p. 47–55.
- 10. Lund BD, Smith JW. A Multi-Stage {CUDA} Kernel for Floyd-Warshall. CoRR [Internet]. 2010;abs/1001.4108. Available from: http://arxiv.org/abs/1001.4108.
- 11. Адуцкевич EB, Лиходед НА, Сикорский AO. К распараллеливанию последовательных программ: распределение массивов между процессорами и структуризация коммуникаций. *Кибернетика и системный анализ*. 2012;48(1):144—163.
 - 12. Воеводин ВВ. Информационная структура алгоритмов. Москва: МГУ; 1997.
- 13. Лиходед НА, Толстиков АА. Формализация коммуникационных операций многомерных циклов. Известия Национальной академии наук Беларуси. Серия физико-математических наук. 2010;3:109—114.

References

- 1. Voevodin VIV, Voevodin VadV. [The fortunate locality of supercomputers]. Otkrytye sistemy [Open systems]. 2013;9:12–15. Russian.
 - 2. Buluc A, Gilberta JR, Budak C. Solving path problems on the GPU. Parallel Computing. 2010;36(5-6):241-253.
- 3. Likhoded NA, Paliashchuk MA. [Estimate of locality of parallel algorithms implemented on GPUs]. *Vestnik Yuzhno-Ural skogo gosudarstvennogo universiteta. Seriya: Vychislitel 'naya matematika i informatika* [Bulletin of the South Ural State University. Series: Computational Mathematics and Software Engineering]. 2016:5(3):96–111. Russian. DOI: 10.14529/cmse160307.
- 4. Likhoded NA, Poleshchuk MA. Method of ranking tiles size parameters of parallel algorithm. *Doklady of the National Academy of Sciences of Belarus*. 2015;59(4):25–33. Russian.
- 5. Kandemir M, Ramanujam J, Irwin M, Narayanan V, Kadayif I, Parikh A. A compiler based approach for dynamically managing scratch-pad memories in embedded systems. *IEEE Transactions on Computer-Aided Design*. 2004;23(2):243–260.

- 6. Baskaran M, Bondhugula U, Krishnamoorthy S, Ramanujam J, Rountev A, Sadayappan P. Automatic data movement and computation mapping for multi-level parallel architectures with explicitly managed memories. In: *Proceedings of the 13th ACM SIGPLAN Symposium on Principles and practice of parallel programming; 2008 February 20–23; Salt Lake City, USA.* New York: ACM; 2008. p. 1–10. DOI: 10.1145/1345206.1345210.
- 7. Poleshchuk MA, Likhoded NA. [Array privatization by computing threads]. In: CSIST'2016. Mezhdunarodnyi kongress po informatike: informatsionnye sistemy i tekhnologii; 24–27 oktyabrya 2016 g.; Minsk, Belarus' [CSIST'2016. International Congress on Computer Science: Information Systems and Technologies; 2016 October 24–27; Minsk, Belarus]. Minsk: Belarusian State University; 2016. p. 883–888. Russian.
- 8. Baskaran M, Ramanujam J, Sadayappan P. Automatic C-to-CUDA code generation for affine programs. In: *Compiler Construction. 19th International Conference. Part of the Joint European Conferences on Theory and Practice of Software; 2010 March 20–28; Paphos, Cyprus.* Berlin, Heidelberg: Springer-Verlag; 2010. DOI: 10.1007/978-3-642-11970-5_14.
- 9. Katz GJ, Kider J. All-pairs shortest-paths for large graphs on the GPU. In: *Proceedings of the 23rd ACM SIGGRAPH/EURO-GRAPHICS symposium on Graphics hardware; 2008 June 20–21; Sarajevo, Bosnia and Herzegovina*. Geneve: Eurographics Association; 2008. p. 47–55.
- 10. Lund BD, Smith JW. A Multi-Stage {CUDA} Kernel for Floyd-Warshall. CoRR [Internet]. 2010;abs/1001.4108. Available from: http://arxiv.org/abs/1001.4108.
- 11. Adutskevich EV, Likhoded NA, Sikorsky AO. [Parallelization of sequential programs: distribution of arrays among processors and structurization of communications]. *Kibernetika i sistemnyi analiz* [Cybernetics and System Analysis]. 2012;48(1):144–163. Russian.
 - 12. Voevodin VV. *Informatsionnaya struktura algoritmov* [Informational structure of algorithms]. Moscow: MGU; 1997. Russian.
- 13. Likhoded NA, Tolstikov AA. Formalization of communication operations of multidimensional loops. *Proceedings of the National Academy of Sciences of Belarus. Physics and Mathematics Series.* 2010;3:109–114. Russian.

Статья поступила в редколлегию 12.06.2018. Received by editorial board 12.06.2018.