

УДК 681.3.06:51

ПАРАЛЛЕЛЬНЫЕ АЛГОРИТМЫ УМНОЖЕНИЯ МАТРИЦ

В.П.Толстьев

В работе обсуждаются способы умножения матриц, которые могут быть эффективно реализованы на многопроцессорных вычислительных системах. Приводятся некоторые сравнительные характеристики рассмотренных алгоритмов.

Запись алгоритмов производится на языке высокого уровня APL [1], удобном для описания операций над векторами, матрицами, деревьями и другими крупными величинами информации.

Скалярные величины будем обозначать строчными латинскими буквами, векторы - подчеркнутыми строчными, а матрицы - подчеркнутыми прописными. Вектор-строка матрицы A обозначается a_i , вектор-столбец - a_j . Размерность вектора a обозначается $v(a)$, размерность матрицы A по строкам - $\mu(A)$, по столбцам - $v(A)$. Все допустимые бинарные операции над векторами и матрицами понимаются как покомпонентные. Дополнительные сведения об операциях и специальных величинах языка APL будут вводиться по мере необходимости. Отметим, что на языке APL операция умножения матриц записывается как частный случай обобщенного матричного произведения и выглядит следующим образом:

$$C \leftarrow A \times B \longleftrightarrow c_{ij} \leftarrow +/(a_i \times b_j),$$

$$i = 1, 2, \dots, \mu(A), j = 1, 2, \dots, v(B), v(A) = \mu(B),$$

где C - результирующая матрица, (i, j) -й элемент которой определяется в результате покомпонентного умножения i -го вектора-строки матрицы A на j -й вектор-столбец матрицы

B и применения к полученному вектору операции р е д у к -
т и в н о г о с л о ж е н и я .

Редуктивная операция выполняется путем последовательного применения ее ко всем компонентам вектора

$$c \leftarrow +/\underline{a} \leftrightarrow c \leftarrow \dots ((\alpha_1 + \alpha_2) + \alpha_3) + \dots + \alpha_v).$$

Редуктивное сложение матрицы по строкам (столбцам) дает в результате вектор, компоненты которого получаютс я р е д у к т и в н ы м сложением строк (столбцов) матрицы:

- 1) редуктивное сложение матрицы по строкам

$$\underline{c} \leftarrow +/\underline{A} \leftrightarrow c_i \leftarrow +/\underline{a}_i, \nu(\underline{c}) = \mu(\underline{A}),$$

- 2) редуктивное сложение матрицы по столбцам

$$\underline{c} \leftarrow +//\underline{A} \leftrightarrow c_j \leftarrow +/\underline{a}_j, \nu(\underline{c}) = \nu(\underline{A}).$$

В данной работе алгоритмы будут записываться в виде последовательности операций языка APL, которая удобна для реализации на многопроцессорной вычислительной системе. В дальнейшем для простоты изложения будем считать, что количество процессоров в системе позволяет производить умножение матриц, не прибегая к разбиению их на подматрицы.

1°. В [2] описан способ умножения матриц, который используется в вычислительной системе ILLIAC IV. Этот способ реализуется на процессорах, настроенных в виде линейной вычислительной системы.

Для нахождения i -й строки результирующей матрицы \underline{C} выбирается i -я строка матрицы \underline{A} , ($1 \leq i \leq \mu(\underline{A})$). Затем поочередно выбирается каждый j -й элемент этой строки ($1 \leq j \leq \nu(\underline{A})$), и с помощью операции р а с п р о с т р а - н е н и я из него формируется вектор-строка размерности $\nu(\underline{B})$, каждая компонента которого равна выбранному элементу. Полученный вектор покомпонентно умножается на соответствующую j -ю строку матрицы \underline{B} ($1 \leq j \leq \nu(\underline{B})$). Образующиеся в результате каждой такой операции векторы покомпонентно складываются между собой. После обработки $\nu(\underline{A})$ -го элемента выбранной строки матрицы \underline{A} образуется редуктивная сумма столбцов матрицы, полученной в результате покомпонентного умножения i -й строки матрицы \underline{A} на все столбцы матрицы \underline{B} , т.е. i -я строка результирующей матрицы \underline{C} . Описанная процедура выполняется до тех пор, пока не будет обработана

последняя строка матрицы A .

Выделение строк и столбцов матриц, а также выделение элементов векторов будем описывать с помощью операции \leftarrow и \rightarrow .

В результате сжатия произвольного вектора A логическим вектором U получается вектор C , содержащий только те компоненты вектора A , которые соответствуют $u_i = 1$. Сжатие записывается как $C \leftarrow A/U$, где $V(U) = V(A)$. Размерность результирующего вектора C равна числу единиц в векторе U .

Сжатие матрицы A логическим вектором U определяется

1) по строкам: $C \leftarrow A/U \rightarrow C_i \leftarrow A_i/u_i$, ($i=1,2,\dots,\mu(A)$), т.е. результирующая матрица C содержит столбцы A , соответствующие $u_j = 1$, а $V(C) = +/U$;

2) по столбцам: $C \leftarrow A/A \rightarrow C_j \leftarrow A_j/a_j$, ($j=1,2,\dots,V(A)$), т.е. результирующая матрица C содержит строки A , соответствующие $u_i = 1$, а $\mu(C) = +/U$.

В качестве сжимающего логического вектора будет использоваться специальный логический вектор $\xi^i(n)$, i -я компонента которого равна единице, а остальные $(n-1)$ компонент — нули.

Операции распространения будем записывать:

- 1) для скаляра по строке $C_i \leftarrow e(n)a_i$,
- 2) для скаляра по столбцу $C_j \leftarrow \phi(n)a_j$,
- 3) для вектора по строкам $C \leftarrow e(n)A_j$,
- 4) для вектора по столбцам $C \leftarrow \phi(n)A_i$,

где n — размерность результирующей величины по соответствующей координате.

С помощью введенных операций алгоритм описанного выше способа умножения матриц можно записать на языке APL следующим образом:

```

<1> i ← 1
<2> j ← 1
от <12> → <3> a_i ← ξi(μ(A))//A
от <10> → <4> b_j ← ξj(μ(B))//B
<5> a_ij ← ξi(V(A))/a_i.

```

$\langle 6 \rangle \underline{d} \leftarrow \Theta(v(\underline{B}))a_{ij}$
 $\langle 7 \rangle \underline{e} \leftarrow \underline{t}_{i.} \times \underline{d}$
 $\langle 8 \rangle \underline{c}_{i.} \leftarrow \underline{c}_{i.} + \underline{e}$
 $\langle 9 \rangle \underline{j} \leftarrow \underline{j} + 1$
 $\kappa \langle 4 \rangle \leftarrow \langle 10 \rangle \underline{j} : v(\underline{A})$
 $\langle 11 \rangle \underline{i} \leftarrow \underline{i} + 1$
 $\kappa \langle 3 \rangle \leftarrow \langle 12 \rangle \underline{i} : \mu(\underline{A})$
 $\langle 13 \rangle$ конец

Результирующая матрица \underline{C} формируется в ячейках $\underline{c}_{i.} \leftarrow \underline{c}_{i.}$.

В рассмотренном способе умножения матриц параллельной обработки подвергаются только векторные величины. При достаточно большом количестве процессоров или незначительных размерах перемножаемых матриц более эффективными могут оказаться операции над матричными величинами.

2°. В работе [3] предлагается структура вычислительной системы, реализующей операции языка APL, в которой предусмотрено специальное устройство для выполнения операций над матрицами. Умножение матриц в этой системе предлагается выполнять способом, аналогичным предыдущему, но используя операции над матричными величинами. Особенность этого метода заключается в том, что выбранная i -я строка матрицы \underline{A} транспонируется и распространяется на $v(\underline{B})$ столбцов. Полученная матрица покомпонентно умножается на матрицу \underline{B} , а результат редуктивно складывается по столбцам, образуя i -ю строку матрицы произведения. Запишем алгоритм этого метода.

$\langle 1 \rangle \underline{i} \leftarrow 1$
 $\rightarrow \langle 2 \rangle \underline{a}_{i.} \leftarrow \underline{t}^i(\mu(\underline{A}))/\underline{A}$
 $\langle 3 \rangle \underline{a}_{.j} \leftarrow \underline{a}_{i.}$
 $\langle 4 \rangle \underline{D} \leftarrow \Theta(v(\underline{B}))\underline{a}_{.j}$
 $\langle 5 \rangle \underline{E} \leftarrow \underline{B} \times \underline{D}$
 $\langle 6 \rangle \underline{c}_{i.} \leftarrow +/\underline{E}$
 $\langle 7 \rangle \underline{i} \leftarrow \underline{i} + 1$
 $\hookleftarrow \langle 8 \rangle \underline{i} : \mu(\underline{A})$
 $\langle 9 \rangle$ конец

3°. С целью экономии памяти вычислительной системы и в тех случаях, когда обращение к памяти в процессе выполнения алгоритма нежелательно, можно применять способ умножения матриц, при котором перемножаемые матрицы целиком находятся в параллельном матричном процессоре в течение всего времени их обработки. В данном случае необходимо произвести некоторые предварительные преобразования исходных величин. Если перемножаются матрицы, которые в общем случае имеют неквадратные формы, то их надо дополнить нулевыми строками или столбцами до квадратных матриц одинаковой размерности ($n \times n$), причем

$$n = \max \{ \nu(A), \mu(A), \nu(B), \mu(B) \}.$$

Преобразования таким образом матрицы будем обозначать A_0 и B_0 . После этого матрица A_0 транспонируется относительно главной диагонали. Покомпонентно умножив транспонированную матрицу A_0 на матрицу B_0 и проведя редуктивное сложение по столбцам, получим элементы главной диагонали результирующей матрицы C . Затем производится циклический сдвиг матрицы B_0 на одну позицию влево по строкам и выполняется следующий шаг алгоритма. Операция циклического сдвига в данном случае заключается в том, что все столбцы матрицы сдвигаются на одну позицию влево, причем крайний левый столбец занимает место крайнего правого. При таком способе выполнения умножения матриц элементы результирующей матрицы будут определяться в следующей последовательности:

1-й шаг - $C_{11} C_{22} C_{33} \dots C_{(n-1)(n-1)} C_{nn}$

2-й шаг - $C_{12} C_{23} C_{34} \dots C_{(n-1)n} C_{n1}$

3-й шаг - $C_{13} C_{24} C_{35} \dots C_{(n-1)1} C_{n2}$

\dots

n -й шаг - $C_{1n} C_{21} C_{32} \dots C_{(n-1)(n-2)} C_{n(n-1)}$

Для того, чтобы после окончания операции элементы результирующей матрицы расположились в необходимом порядке, можно полученный на каждом шаге вектор элементов записывать в главную диагональ матрицы, в которой формируется результат, и после выполнения очередного шага производить циклический сдвиг содержимого этой матрицы на одну позицию влево по строкам.

Для описания рассматриваемого способа умножения матриц необходимы специальные векторы $\underline{E}(n)$ - полный единичный

вектор размерности n , каждая компонента которого равна единице, и $\underline{L}(n)$ — интервальный вектор, представляющий собой отрезок натурального ряда чисел от 0 до $(n-1)$. Операция транспонирования матрицы относительно главной диагонали обозначается $\underline{C} \leftarrow \underline{A}$. Операции циклического сдвига матриц записываются:

- 1) по строкам вправо — $\underline{C} \leftarrow \underline{k} \downarrow \underline{A}$,
- 2) по строкам влево — $\underline{C} \leftarrow \underline{k} \uparrow \underline{A}$,
- 3) по столбцам вниз — $\underline{C} \leftarrow \underline{k} \downarrow \underline{A}$,
- 4) по столбцам вверх — $\underline{C} \leftarrow \underline{k} \uparrow \underline{A}$,

где \underline{k} — целочисленный вектор, каждая компонента которого указывает, на какое число позиций сдвигается соответствующая строка (столбец) матрицы.

Алгоритмы данного способа умножения матриц можно записать в виде такой последовательности операций:

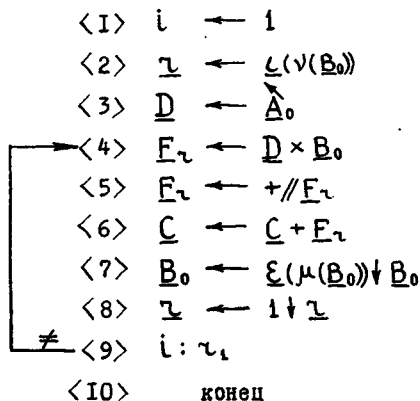
```

<1> i ← 1
<2> D ← A0
<3> E ← D × B0
<4> n1 ← +//E
<5> N ← Et(v(B0)) × n1
<6> N ← Lo(v(B0)) ↓ N
<7> C ← C + N
<8> C ← E(μ(A0)) ↑ C
<9> B0 ← E(v(B0)) ↑ B0
<10> i ← i + 1
<11> i : μ(A0)
<12> . конец

```

Операции <5> + <8> в приведенном алгоритме можно устранить, если снабдить параллельный матричный процессор индекс-регистром, состоящим из соединенной в кольцо цепочки ячеек, каждая из которых управляет соответствующим столбцом процессоров и может хранить числа от 1 до $v(B_0)$. Число, хранящееся в ячейке индекс-регистра, показывает, в каком процессоре соответствующего столбца должен формироваться результат редуktiv-

ного сложения на данном шаге алгоритма. Перед началом операции в индекс-регистр заносится интервальный вектор $\underline{L}(\underline{V}(\underline{B}))$, начинающийся с единицы. Поэтому на первом шаге алгоритма значения элементов результирующей матрицы формируются в процессорах, находящихся на главной диагонали. В таком же порядке эти элементы переписываются в матрицу результата с помощью покомпонентной операции сложения. Затем производится кольцевой сдвиг содержимого индекс-регистра на одну позицию влево и выполняется следующий шаг алгоритма. Элементы результирующей матрицы в соответствии с показаниями индекс-регистра будут занимать правильные места в матрице результата. Такая модификация позволяет существенно упростить алгоритм рассматриваемого способа умножения матриц. Теперь он будет выглядеть следующим образом:



В приведенной записи алгоритма \underline{L} означает индекс-регистр, F_r - управляемый матричный процессор.

4°. Рассмотрим еще один способ умножения матриц, который заключается в последовательном накоплении значений элементов результирующей матрицы. Пусть $\underline{V}(\underline{A}) = \underline{\mu}(\underline{B}) = n$. Тогда каждый элемент результирующей матрицы находится как сумма произведений соответствующих пар элементов перемножаемых матриц:

$$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj} \quad (1)$$

На каждом шаге алгоритма этого способа умножения матриц определяются очередные слагаемые для всех элементов результирующей

матрицы, которые складываются с соответствующими ранее накопленными суммами. Для того, чтобы найти все слагаемые k -го шага, выбираются k -й столбец матрицы A и k -я строка матрицы B . Первый из этих векторов распространяется на $\nu(B)$ позиций по строкам, а второй на $\mu(A)$ позиций по столбцам. Полученные матрицы покомпонентно перемножаются, в результате чего образуется матрица k -х членов для всех сумм вида (1). Алгоритм этого способа можно записать в виде последовательности операций:

```

<1> i ← 1
      |
      |→ <2>  $a_{.j} \leftarrow \sum^i (\nu(A))/A$ 
      |
      |→ <3>  $D \leftarrow \Theta(\nu(B)) a_{.j}$ 
      |
      |→ <4>  $b_{i.} \leftarrow \sum^i (\mu(B))/B$ 
      |
      |→ <5>  $E \leftarrow \Phi(\mu(A)) b_{i.}$ 
      |
      |→ <6>  $G \leftarrow D \times E$ 
      |
      |→ <7>  $C \leftarrow C + G$ 
      |
      |→ <8>  $i \leftarrow i + 1$ 
      |
      |↙ <9>  $i : \nu(A)$ 
      |
      |→ <10> конец

```

В заключение сравним по быстродействию рассмотренные способы умножения матриц. Без большого ущерба для точности оценок будем считать, что перемножаются квадратные матрицы размерности $(n \times n)$. Такую же размерность должны иметь матричные параллельные процессоры. Чтобы сравнивать все алгоритмы при одинаковом количестве оборудования, будем рассматривать линейную вычислительную систему как n линейных групп по n процессоров в каждой, причем на каждой группе процессоров вычисляется одна строка результирующей матрицы. При такой организации линейной вычислительной системы для умножения матриц на каждой группе процессоров обрабатывается одна из строк матрицы A и все строки матрицы B , т.е. внешний цикл алгоритма 1^0 по i распараллеливается, в результате чего этот способ умножения матриц по существу сводится к способу 4^0 .

В расчетах можно пренебречь временем выполнения вспомогательных операций таких, как распространение, сдвиг, выделение и других потому, что а) количество таких операций во всех

рассмотренных алгоритмах приблизительно равно, б) это очень быстрые операции, в) многие из них могут быть совмещены во времени с основными операциями — умножением и делением.

При введенных ограничениях алгоритмы 1^0 и 4^0 будут содержать, без учета вспомогательных операций, n групповых операций умножения и n групповых операций сложения. Каждая из таких операций состоит из n^2 одноименных элементарных операций, выполняемых одновременно. Обозначим время выполнения одной групповой операции умножения и сложения соответственно через $t(\times)$ и $t(+)$. Тогда полное время выполнения алгоритмов 1^0 и 4^0 будет равно:

$$T_{1,4} = n \cdot t(\times) + n \cdot t(+). \quad (2)$$

Алгоритмы 2^0 и 3^0 также требуют выполнения n аналогичных групповых операций умножения, поэтому время выполнения этих операций во всех алгоритмах одинаково и равно $n \cdot t(\times)$. Для выполнения операций сложения в алгоритмах 2^0 и 3^0 требуется больше времени, чем в оценке (2). Это вызвано тем, что на каждом шаге алгоритмов 2^0 и 3^0 необходимо выполнить непосредственно редуктивное сложение векторов, которое невозможно свести к одной групповой операции сложения. Эту операцию можно максимально распараллелить путем последовательного разбиения элементов редуктивно складываемых векторов и промежуточных сумм на пары слагаемых. Поскольку после каждого проведенного таким образом сложения число пар слагаемых уменьшается вдвое, то при размерности векторов, равной n , вся операция редуктивного сложения векторов потребует $\log_2 n^*$ групповых сложений, где n^* — ближайшая большая или равная n степень 2. Операцию редуктивного сложения векторов необходимо выполнять на каждом шаге алгоритмов 2^0 и 3^0 . Следовательно, полное время выполнения этих алгоритмов равно:

$$T_{2,3} = n \cdot t(\times) + n(\log_2 n^*) t(+). \quad (3)$$

Исходя из полученных оценок, можно сделать вывод, что наибольшее быстродействие достигается с помощью алгоритмов 1^0 и 4^0 . Отметим, что в действительности время выполнения алгоритма 3^0 будет несколько больше приведенной оценки (3), поскольку этот алгоритм требует значительных преобразований исходных матриц перед их умножением этим способом.

Наиболее серьезные практические трудности связаны с обработкой матриц, размерность которых не совпадает с размерностью параллельных процессоров. В этом случае встают такие вопросы, как простота разбиений исходных матриц, удобство обращения к промежуточным результатам, оптимизация загрузки процессоров и т.д. Поэтому основным требованием к алгоритмам становится их гибкость при решении перечисленных вопросов. Из рассмотренных нами алгоритмов несомненными преимуществами в этом отношении обладает первый алгоритм, т.к. линейная вычислительная система легко перестраивается на различное число групп процессоров, позволяя тем самым устранить простой оборудования. Кроме того, если матричные параллельные процессоры позволяют перемножать без разбиений матрицы размерностью $(n \times n)$, то линейная вычислительная система на таком же количестве процессоров позволяет без разбиений получить произведение матриц

$$C(m \times n^2) \leftarrow A(m \times \ell) \times B(\ell \times n^2),$$

где размерности ℓ и m ограничиваются лишь длиной машинного слова и памятью процессоров. При этом надо учесть, что метод разбиения матриц на подматрицы в случае алгоритма 1^0 наиболее прост. Однако, в тех случаях, когда операция умножения матриц является лишь одной из многочисленных операций алгоритма решения более сложной задачи, целесообразным может оказаться применение алгоритмов 2^0 или 3^0 , т.к. каждый шаг этих алгоритмов дает окончательные значения элементов результирующей матрицы, которые можно сразу использовать для дальнейших вычислений.

Л и т е р а т у р а

1. Iverson K.E. A programming language. New York, Wiley, 1962.
2. McIntyre D.E. An introduction to the ILLIAC-IV computer. - Datamation, 1970, 16, N4, p. 60-67.
3. Thurber K.J., Myrna J.W. System design of a cellular APL computer. - IEEE Trans. on Comp., 1970, C-19, N4, p. 291-301.

Поступила в ред.-изд.отдел
3 апреля 1972 г.