

А.В. Чистяков, И.С.Ислямова

МЕТОД И ТЕХНОЛОГИИ ПАРАЛЛЕЛЬНОГО ПРОГРАММИРОВАНИЯ ПРИ РЕШЕНИИ ПРИКЛАДНЫХ ЗАДАЧ

**Национальный
авиационный университет**

**кафедра инженерии про-
граммного обеспечения**

**Научный руководитель
Иванова Л.Н.
(к.т.н., доцент)**

Постановка задачи

На сегодняшний день ни одна отрасль народного хозяйства, а также любого производства не может эффективно работать без использования компьютерной техники, а тем более решать конструкторские задачи, задачи моделирования и создания новых образцов. Даже, на первый взгляд, несложные задачи, которые встречаются в повседневной жизни, не всегда могут быть решены эффективно классическими методами и средствами линейного или структурного программирования. Также и наличие мощной вычислительной техники не гарантирует быстрого выполнения поставленной задачи средствами среды разработки программного обеспечения (ПО). В качестве примера можно привести приложение, написанное на языке программирования C++ и скомпилированное в режиме многопоточного выполнения, однако общее время исполнения сократится незначительно. Поэтому возникает задача максимально эффективного использования вычислительных ресурсов ЭВМ как при разработке нового ПО, так и при модернизации существующего. Одним из решений этой задачи может быть использование технологий и методов параллельного программирования.

Анализ существующих топологий многоядерных компьютеров

На сегодняшний день в мире существует множество классов и типов компьютеров, их все можно классифицировать по количеству потоков команд и по количеству потоков данных, которые обрабатывает одновременно система (классификация Флинна):

- ОКОД — Вычислительная система с одиночным потоком команд и одиночным потоком данных (SISD, Single Instruction stream over a Single Data stream).
- ОКМД — Вычислительная система с одиночным потоком команд и множественным потоком данных (SIMD, Single Instruction, Multiple Data).
- МКОД — Вычислительная система со множественным потоком команд и одиночным потоком данных (MISD, Multiple Instruction Single Data).
- МКМД — Вычислительная система со множественным потоком команд и множественным потоком данных (MIMD, Multiple Instruction Multiple Data).

Самыми распространенными и эффективными промышленными системами являются MIMD-компьютеры. В состав такого компьютера входит несколько процессоров, которые функционируют асинхронно и независимо друг от друга. В любой момент времени различные процессоры могут выполнять различные команды над разными частями одних и тех же данных. MIMD-компьютеры могут быть как однородные, состоящие из одинаковых узлов, так и разнородные (полисистемы), состоящие из различных по аппаратной составляющей узлов. Разнородные системы часто называют Beowulf-системами или Beowulf-кластерами. Beowulf-кластер состоит из широко распространенного аппаратного обеспечения, работающие под управлением операционной системы с открытым исходным кодом (например, GNU/Linux или FreeBSD). Основным до-

стоинством таких систем является их дешевизна, высокая производительность и возможность комбинировать любые модели компьютеров в один кластер.

MIMD-компьютеры относятся к компьютерам с разделенной памятью. А персональные компьютеры и малые сервера, привычные обычному пользователю, относятся к SIMD-системам с общей памятью.

В последние годы очень активно развиваются гибридные системы и вычисления на графических ускорителях (видеокартах). Гибридная система представляет собой персональный компьютер (ПК), сервер или кластер, на котором установлены специализированные программно-аппаратные комплексы, позволяющие выполнять задачи не только на центральном процессоре (CPU), но и на процессоре видеокарты (graphics processing unit, GPU). В этом случае работу вычислительной системы можно организовать таким образом, что бы поставленная задача могла выполняться параллельно как на отдельно взятом GPU, так и совместно с CPU. Совместное использование CPU с GPU позволяет повысить эффективность работы ПО в десятки раз. Это связано с тем, что процессор видеокарты имеет в своем составе большое количество арифметико-логических устройств, которые специализированы под обработку больших массивов данных. Но при решении прикладных задач (моделирования физики процессов, генома, фармацевтики, погоды, и т.д.) не всегда достаточно даже суммарной производительности CPU и GPU. Если установить в узлы кластера видеокарты и использовать несколько GPU параллельно, то можно повысить эффективность решения поставленной задачи на несколько порядков.

Графические процессоры можно рассматривать как мощные параллельные SIMD-процессоры (Single Instruction Multiple Data), способные выполнять одну и ту же операцию одновременно над несколькими значениями однородных данных. То есть SIMD-процессор получает на вход поток однородных данных и параллельно обрабатывает их, порождая тем самым выходной поток.

Анализ существующих технологий и средств параллельного программирования

На сегодняшний день можно выделить несколько основных технологий параллельного программирования для систем с разделенной памятью (Message Passing Interface, Parallel Virtual Machine), для систем с общей памятью

(Open Multi-Processing), а так же для программирования видеокарт и GPU (CUDA, ATI Stream Technology).

Message Passing Interface (MPI, интерфейс передачи сообщений) [1] — программный интерфейс для передачи информации, который позволяет обмениваться сообщениями между процессорами, выполняющими одну задачу. Разработан Уильямом Гроуппом, Эвином Ласком и другими.

MPI является наиболее распространённым стандартом интерфейса обмена данными в параллельном программировании. Существуют его реализации для большого числа компьютерных платформ и для языков программирования

Фортран 77/90, Си и Си++. Основным средством коммуникации между процессами в MPI является передача сообщений друг другу. Стандартизацией MPI занимается MPI Forum. В стандарте MPI описан интерфейс передачи сообщений, который должен поддерживаться как системой, на которой выполняется ПО, так и самим ПО пользователя. В настоящее время существует большое количество бесплатных и коммерческих дистрибутивов MPI: MPICH, LAM, HPVM, OpenMPI, WMPi и другие.

Одной из основных особенностей MPI является то, что эта система лучше всего работает в однородных системах. Это связано с тем, что используя базовый пакет MPI, система будет работать до тех пор, пока задача не будет решена на самом маломощном узле разнородного кластера. Вследствие чего большая часть вычислительного ресурса будет простаивать, что не целесообразно. Средствами MPI решить задачу балансировки вычислительной нагрузки сложно, для этого требуется писать специальные алгоритмы решения этой задачи.

Parallel Virtual Machine [2] (PVM, параллельная виртуальная машина) — является основой вычислительной среды Beowulf-кластера, который представляет собой пакет программ и позволяет использовать связанный в локальную сеть набор разнородных компьютеров, работающих под операционной системой Unix, как один большой параллельный компьютер. Таким образом, проблема больших вычислений может быть весьма эффективно решена за счет использования совокупной мощности и памяти большого числа компьютеров. Пакет программ PVM легко переносится на любую платформу, начиная от laptop и до CRAY.

PVM можна визначити як частину засобів реального висувального комплексу (процесори, пам'ять, периферійні пристрої і т.д.), призначену для виконання множини завдань, учасників в отриманні загального результату висувальних. В загальному випадку кількість завдань може перевищувати кількість процесорів, включених в PVM. Крім того, в склад PVM можна включати досить різноманітні висувальні машини, несумісні за архітектурі даних. Інакше кажучи, паралельною віртуальною машиною може стати як окремо взятий персональний комп'ютер (ПК), так і локальна мережа, включаючи в себе суперкомп'ютери з паралельною архітектурою, універсальні ЕВМ, графічні робочі станції і всі ті ж маломощні ПК. Важливо лише, щоб у включених в PVM висувальних засобах була інформація в використовуваному програмному забезпеченні PVM. Завдяки цьому програмному забезпеченню користувач може вважати, що він спілкується з однією висувальною машиною, в якій можливо паралельне виконання множини завдань.

Використання PVM більш бажано в кластерах типу Beowulf. Це пов'язано з тим, що при розробці програми програміст повинен враховувати апаратну особливість кожного вузла кластера, а це не завжди зручно. Але найбільш вирішується завдання з простим висувальним потужностям.

Для програмування в системах з загальною пам'яттю використовується технологія OpenMP.

OpenMP (Open Multi-Processing) [3,4] - відкритий стандарт для розпаралелювання програм, написаних на мовах програмування Си, Си++ і Фортран. Описує набір директив компілятора, бібліотечних процедур і змінних оточення, які призначені для програмування багатопотокових програм на багатопроцесорних системах з загальною пам'яттю. Розробку специфікації OpenMP ведуть кілька великих виробників висувальної техніки і ПО, чья робота регулюється некомерційною організацією, званою OpenMP Architecture Review Board (ARB).

OpenMP реалізує паралельні висувальні з допомогою багатопотоковості, в якій «головний» (master) потік створює набір підлеглих (slave) потоків і завдання розподіляються між ними. Припускається, що потоки виконуються паралельно на машині з кількома процесорами (кількість процесорів

не обов'язково повинно бути більше або рівно кількості потоків).

MPI і OpenMP технології можуть бути використані разом для підвищення ефективності функціонування багатоядерних вузлів в кластерах, так як MPI орієнтована на системи з розділеною пам'яттю, тобто коли витрати на передачу даних великі, а OpenMP орієнтована на системи з загальною пам'яттю (багатоядерні ПК).

Завдання, виконувані потоками паралельно, також як і дані, потрібні для виконання цих завдань, описуються з допомогою спеціальних директив препроцесора відповідного мови — прагм.

Ключовими елементами OpenMP є:

- конструкції для створення потоків (директива `parallel`);
- конструкції розподілу роботи між потоками (директиви `DO/for` і `section`);
- конструкції для управління роботою з даними (вираження `shared` і `private`);
- конструкції для синхронізації потоків (директиви `critical`, `atomic` і `barrier`);
- процедури бібліотеки підтримки часу виконання (наприклад, `omp_get_thread_num`);
- змінні оточення (наприклад, `OMP_NUM_THREADS`).

OpenMP - це ідеальний засіб для модернізації існуючого ПО, що функціонує в однопроцесорних системах. Так як ця технологія використовує прагми для управління потоками і тому немає потреби масштабно редагувати вихідний код.

CUDA (Compute Unified Device Architecture) — це технологія компанії NVIDIA, заснована на розширенні мови Си і яка дає можливість управління набором інструкцій і пам'яттю графічного прискорювача для організації паралельних висувальних. CUDA може бути використана на графічних процесорах відеокарт (відеоприскорювачів) GeForce восьмого покоління і старші (серії GeForce 8, GeForce 9, GeForce 200), а також Quadro і Tesla.

Трудоємкість програмування GPU з допомогою CUDA досить велика, однак вона нижче, ніж з ранніми GPGPU (General-purpose graphics processing units — «GPU загального призначення») рішеннями. При створенні ПО з використанням CUDA потрібно враховувати його виконання на кількох мультипроцесорах, так же як і при MPI про-

граммуванні, але без розділення даних, які зберігаються в загальній відеопам'яті. І так як CUDA програмування для кожного мультипроцесора подібно OpenMP програмуванню, воно потребує хорошого розуміння організації пам'яті. Але, звичайно ж, складність розробки нової та переносу існуючого ПО на CUDA в великій ступені залежить від розв'язуваної задачі.

Технологія програмування CUDA використовується для організації масивно-паралельних операцій на GPU [5]. При цьому послідовна частина програми виконується на CPU, а масивно-паралельні обчислення організуються на GPU, як набір ниток, одночасно виконуваних потоків (threads). При цьому кожній нитці відповідає один елемент обчислюваних даних. Нитки об'єднуються в сітки (grid) і блоки (block). Кожний блок може бути одномерним, двимірним або тримірним. Таким чином, задача в CUDA розбивається на декілька окремих підзадач. Кожній такій підзадачі відповідає свій блок ниток, і взаємодія між нитками може відбуватися тільки в межах одного блоку на швидкій розподільчій пам'яті.

Технологію CUDA може використовувати будь-який програміст, який знає мову Си. Прийдеться тільки звикнути до іншої парадигми програмування, властивій паралельним обчисленням. Але якщо алгоритм за принципом добре розпаралелюється, то вивчення та витрати часу на програмування на CUDA повернуться в декілька разів.

Технологія CUDA включає в себе і дозволяє наступне:

- уніфікований програмно-апаратний комплекс для паралельних обчислень на відеочіпах NVIDIA;
- набір інструментів, які підтримують роботу з GPU різних типів пристроїв;
- підтримує програмування на мові програмування Си;
- стандартні бібліотеки чисельного аналізу FFT (швидке перетворення Фур'є) і BLAS (лінійна алгебра);
- оптимізований обмін даними між CPU і GPU;
- взаємодія з графічними API (application programming interface, Інтерфейс прикладного програмування) OpenGL і DirectX;

- підтримку 32- і 64-бітних операційних систем: Windows XP, Windows Vista, Linux і MacOS X;

- можливість розробки ПО на низькому рівні.

В доповнення до властивості підтримки операційних систем необхідно уточнити, що офіційно підтримуються всі основні дистрибутиви Linux (Red Hat Enterprise Linux 3.x/4.x/5.x, SUSE Linux 10.x), але CUDA прекрасно працює і на інших збірках операційних систем (ОС): Fedora Core, Ubuntu, Gentoo і др.

Середовище розробки CUDA (CUDA Toolkit) включає:

- компілятор nvcc;
- бібліотеки FFT і BLAS;
- профілювальник;
- відладчик gdb для GPU;
- драйвер CUDA runtime в комплекті стандартних драйверів NVIDIA;
- посібник з програмування;
- CUDA Software Development Kit (SDK, комплект засобів розробки ПО) (вихідний код, утиліти і документація).

В складі CUDA SDK входять приклади: паралельна бітонна сортування (bitonic sort), транспонування матриць, паралельне префіксне суммування великих масивів, свертка зображень, дискретне вейвлет-перетворення, приклад взаємодії з OpenGL і Direct3D, використання бібліотек CUBLAS і CUFFT, обчислення ціни опціону (формула Блэка-Шоулза, біноміальна модель, метод Монте-Карло), паралельний генератор випадкових чисел Mersenne Twister, обчислення гистограм великого масиву, шумоподавлення, фільтр Собеля (знаходження меж).

Технологія FireStream - це програмно-апаратна обчислювальна архітектура компанії AMD. AMD FireStream (раніше ATI FireStream і AMD Stream Processor) - є потіковим процесором, розробленим компанією ATI [6,7].

Областями застосування FireStream є задачі, які потребують великого обчислювального ресурсу, такі як фінансовий аналіз або обробка сейсмічних даних. Використання потікового процесора дозволило збільшити швидкість деяких фінансових розрахунків в 55 разів порівняно з рішенням тієї ж задачі силами тільки центрального процесора.

К системе программирования ATI Stream входит набор приложений и процессоры ATI Stream. На рис. 1 представлено взаимодействие всех компонентов системы программирования ATI Stream. Система программирования ATI Stream обеспечивает конечных пользователей и разработчиков гибкими средствами для пользо-

вания графическими процессорами. ПО ATI поддерживает открытые системы и открытые стандарты. Таким образом, открытая стратегия ATI дает возможность разработчикам внедрять программное обеспечение, разработанное для ATI Stream по лицензии GPL.

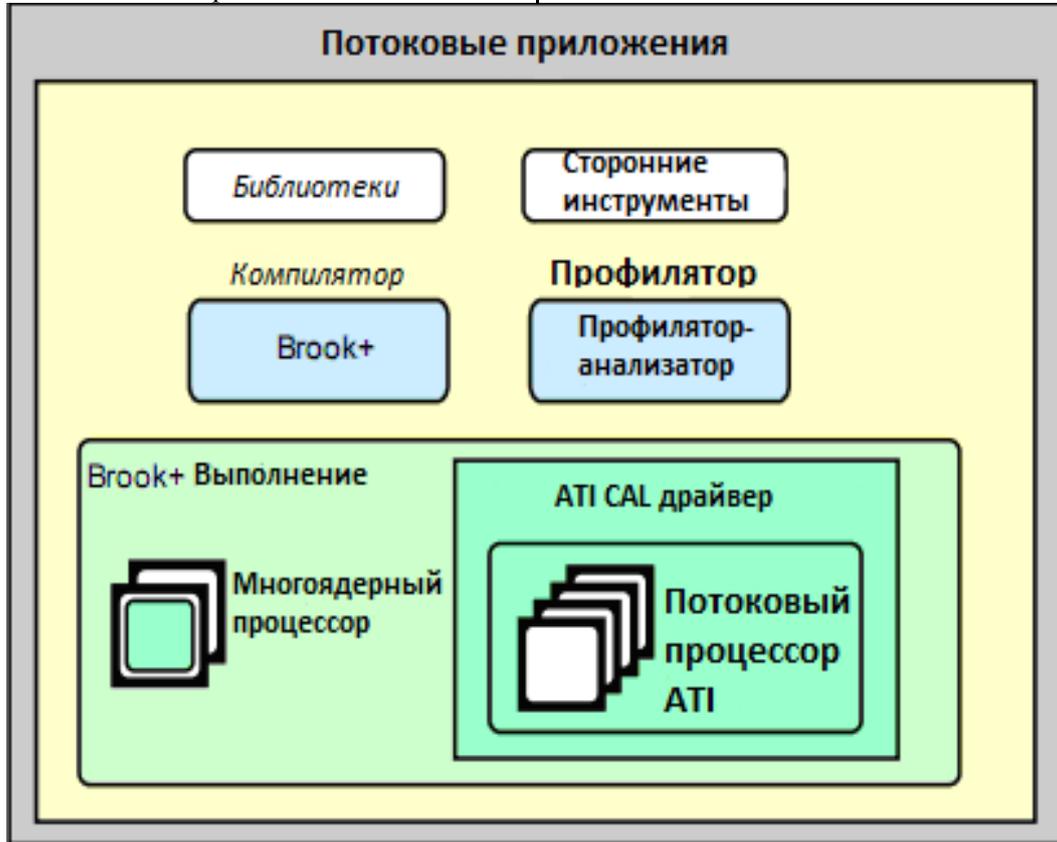


Рис. 1 Взаимодействие систем компонентов ATI Stream

В состав программного обеспечения входит:

- компилятор brook + с дополнениями для ATI устройств;
- драйвер устройства для ATI процессоров (ATI Compute Abstraction Layer – CAL);
- профилятор-анализатор быстрого действия - Stream Kernel Analyzer;
- вычислительная библиотека - AMD Core Math Library (ACML);

Программирование графических процессоров ATI Stream проводится с помощью технологии программной модели шейдеров. Поточковые процессоры выполняют пользовательские приложения, называемые потоковым ядром (stream kernel). Поточковые процессоры могут использоваться для неграфических вычислений, используя модель программирования SIMD. В этой модели программирования, которая называется потоковым программиро-

ванием, массивы входных данных, хранящихся в общей памяти, разделяются между определенным количеством SIMD процессоров, которые в свою очередь выполняют потоковые ядра и записывают результат в общую память.

Каждый экземпляр ядра в текущем процессоре, называется нитью (thread). Существует определенный квадратный регион, в котором отражаются все нити, который в свою очередь называется домен исполнения (domain of execution).

Потоковый процессор формирует очередь нитей для определенной группы ниточных процессоров (thread processors) до того, пока применение не закончит работу. Упрощенную модель вычислений ATI Stream можно увидеть на рис. 2.

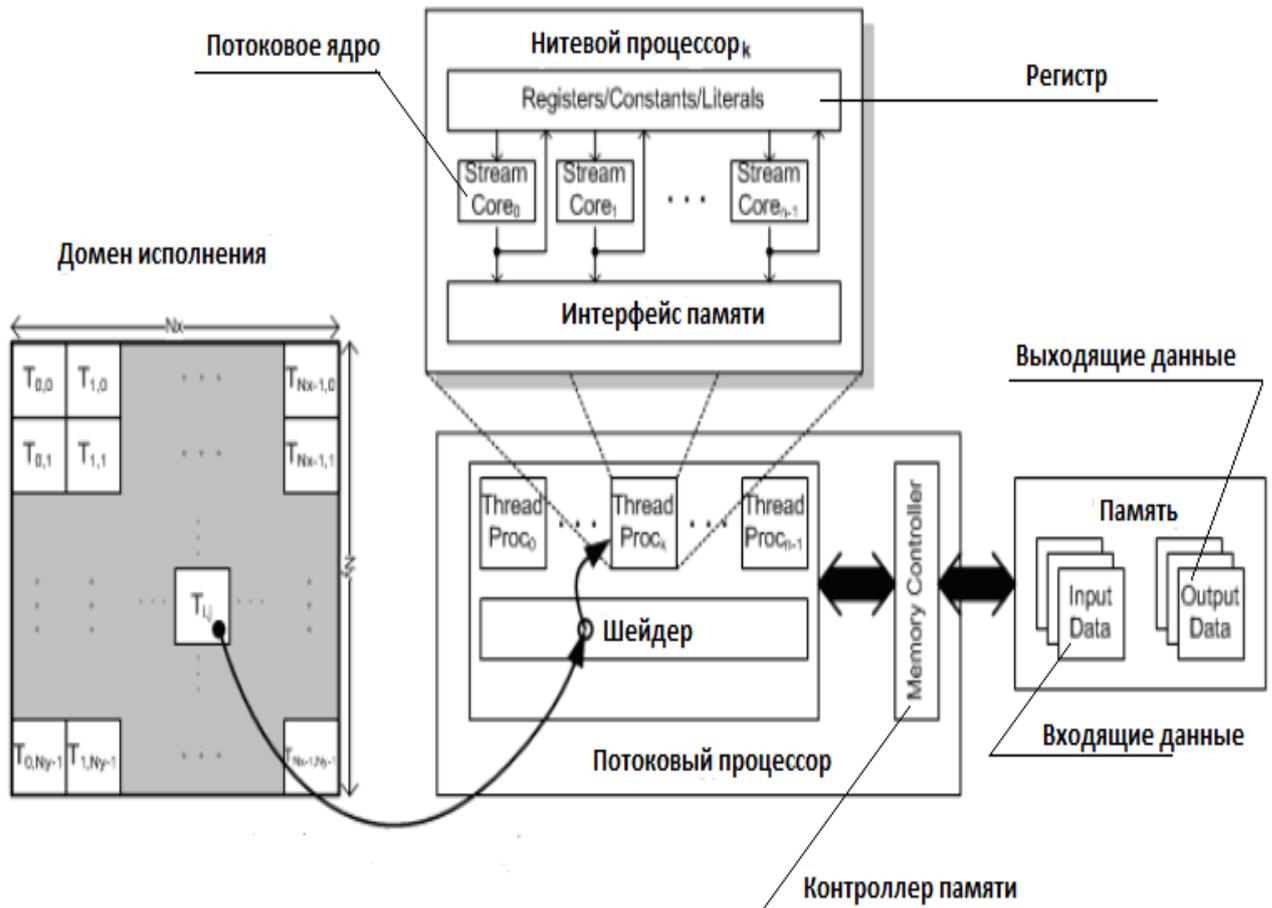


Рис. 2 Упрощена модель вичислень ATI Stream

Открытый параллельный компилятор Brook +. Brook + дает возможность явного распараллеливания, используя расширения языка ANSI C. Модель вычислений Brook +, так называемая потоковая модель, разрабатывается параллельно с традиционными методами параллельного программирования и позволяет:

- организовать параллелизм по данным., таким образом позволяет использовать свойство SIMD-архитектуры.
- использование повышенной арифметической точности, таким образом позволяет разработку компьютерных алгоритмов с максимальной точностью результата и сокращением времени коммуникации между ядрами.

В языке программирования Brook + существует два ключевых элемента:

- Поток.
- Ядро.

Пример:

```
kernel void sum (float <> a, float <> b, float <> c)
{
```

```
c = a + b
}
```

Такой пример кода, написанный на языке Brook +, добавляет два потока и записывает их в выходной поток.

На рис. 3 изображены элементы языка программирования Brook +, в который входят [7]:

- brcc - компилятор программ Brook + с расширением. br, что превращает код Си в код, понятный графическим процессорам (у этого компилятора есть функция создания виртуальной среды, позволяющей отлаживать программы);

brt - библиотека, которая дает возможность пользоваться функциями для работы с графическими процессорами.

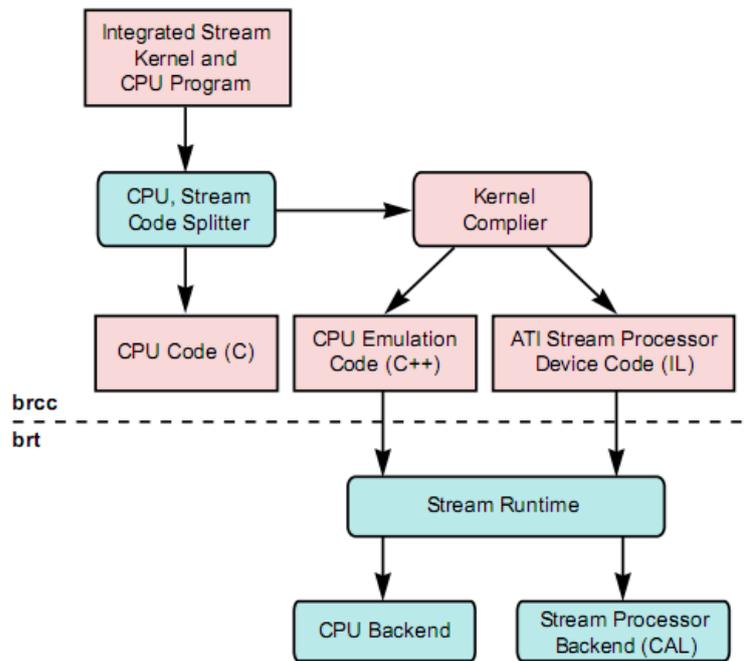


Рис. 3 Елементи языка программирования Brook+

Абстрактный уровень вычислений (ATI Computation Abstraction Layer – CAL).

Как показано на рис. 4 существует абстрактный уровень вычислений, позволяющий пользоваться функциями для управления ядра-

ми. Такой абстрактный уровень представлен драйверами и библиотекой. Таким образом, такая модель позволяет пользователю управлять процессом выполнения программы и тем самым повысить ее производительность.

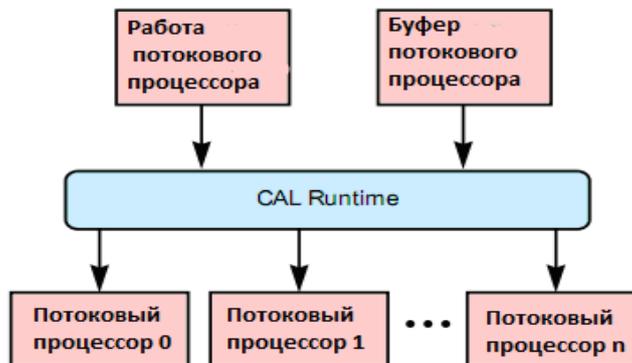


Рис. 4 Модель коммуникации среды и процессоров

CAL позволяет:

- генерировать инструкции для управления аппаратным обеспечением;
- управлять аппаратным обеспечением;
- управлять ресурсами;
- выполнять инструкции;
- поддерживать работу многих устройств;
- интероперабельность с разными 3D платформами.

К CAL входят функции языка Си и наборы типов данных, позволяющих высокоуровневым приложениям управлять процессами и количеством памяти на устройстве, где выполняется параллельное приложение.

AMD Core Math Library (ACML). В состав среды разработки входит математическая библиотека линейной алгебры ACML, которая имеет функции BLAS. Такая библиотека вклю-

чаєт основні математическі функції. Она оптимізована для платформи АТІ.

В склад цієї бібліотеки входять наступні модулі:

- весь набір функцій BLAS;
- набір функцій LAPACK;
- функції швидкого перетворення Фур'є;
- функції генератора випадкових чисел.

Потокові процесори.

На рис. 5 показана спрощена структура поточкового процесора. Існує велика кількість поточкових процесорів, але більшість з них мають наступні модулі:

- SIMD двигателі;

- нитеві процесори;
- Поточкові ядра.

В склад поточкового процесора входять групи SIMD двигателів, які в свою чергу мають визначене кількість нитевих процесорів, відповідальних за виконання ядер / модулів, кожне в окремому потоці. Ниточний процесор включає велику кількість поточкових ядер, відповідальних за розрахунок різної точності. Кожний поточковий процесор виконує одну інструкцію, але з різним набором даних, який показує властиву їм SIMD-архітектуру.

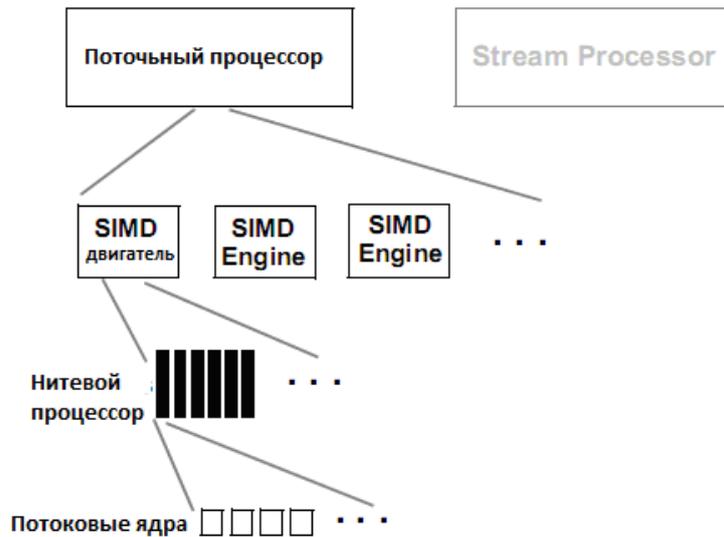


Рис. 5 Спрощена структура поточкового процесора

Існують деякі рекомендації (правила) по переносу ПО на багатопотокові технології, наприклад, планування циклу і відповідно баланс навантаження. Ці правила однаково повинні застосовуватися для будь-якої з вище наведених технологій.

Метод інкрементного програмування і правила паралельного програмування

Створення ефективного алгоритму – найскладніша частина створення будь-якого ПО, особливо паралельного. Тому для спрощення написання і модернізації вже існуючого ПО часто використовують метод інкрементного програмування. Інкрементне програмування – це процес детального розбору алгоритму або ПО на складові частини таким чином, щоб отримати незалежні один від одного блоки даних і/або команд. По результатам аналізу ПО програміст дописує необхідні блоки коду, відповідальні за розпаралелювання, керуючись правилами, при-

веденими нижче. Для простоти викладу, будемо вважати, що ми модернізуємо деяке ПО за допомогою технології OpenMP [8].

Першим правилом паралельного програмування є правило про баланс навантаження на вузли системи. Баланс навантаження (розподілення робочого навантаження порівну між потоками), є одним з найбільш важливих атрибутів паралельного виконання програми. Це правило має велике значення, оскільки гарантує роботу всіх процесорів більшу частину часу. Без балансу навантаження деякі потоки можуть завершити роботу значно раніше інших, що призводить до простою обчислювальних ресурсів і втрати продуктивності.

В циклах відсутність балансу навантаження зазвичай є наслідком різниці часу обчислення в різних ітераціях циклу. Розброс часу обчислення в ітераціях циклу зазвичай легко визначити, вивчив вихідний код. В більшості випадків ми побачимо, що

итерации цикла занимают одинаковое время. Если это не так, то можно найти наборы итераций, занимающие одинаковое время. Например, иногда набор всех четных итераций занимает примерно столько же времени, как и набор всех нечетных итераций. Аналогично, первая половина итераций цикла может занимать примерно столько же времени, как и вторая половина. С другой стороны, может оказаться невозможным найти наборы итераций, имеющие одинаковое время выполнения. Независимо от того, какой из этих случаев имеет место в конкретном приложении, необходимо предоставить OpenMP эту дополнительную информацию о планировании цикла, чтобы он мог правильно распределить итерации цикла между потоками (и, следовательно, между процессорами) для оптимизации распределения нагрузки.

По умолчанию, OpenMP предполагает, что все итерации цикла занимают одинаковое время. В результате OpenMP распределяет итерации цикла между потоками примерно поровну таким образом, чтобы минимизировать вероятность возникновения конфликтов памяти вследствие ее неправильного совместного использования. Это возможно, поскольку итерации цикла обычно обращаются к памяти последовательно. Поэтому при разделении цикла на две большие части (например, на первую и вторую половины) при использовании двух потоков вероятность наложения памяти оказывается наименьшей. Однако, хотя это и может быть наилучшим вариантом во избежание конфликтов памяти, с точки зрения баланса нагрузки это может быть плохим выбором. К сожалению, обратное тоже справедливо. То, что хорошо для баланса нагрузки, может быть плохо для работы с памятью. Поэтому инженерам по производительности необходимо найти баланс между оптимальным использованием памяти и оптимальным распределением нагрузки, измеряя производительность, чтобы определить, какие методы дают наилучшие результаты.

```
void m_mult_m (double *a, double *b, double *c, int n)
```

```
{
    int i, j, k;
    for(i = 0; i < n; i++)
        for(j = 0; j < n; j++)
            c[i][j] = 0;
    for (i = 0; i < n; i++){
        for (j = 0; j < n; j++){
            double s = 0;
            for (k = 0; k < n; k++)
                s += a[i][k] * b[k][j];
        }
    }
}
```

Как можно было понять из изложенного выше, вторым правилом параллельного программирования является то, что программу или функцию невозможно эффективно распараллелить если она обращается к общим ресурсам. Например, если нам надо чтобы несколько потоков в процессе решения одной задачи выводили на консоль, допустим, контрольные суммы в определенном порядке. Так как управление консолью будет передаваться каждому потоку в порядке очереди, которая формируется по мере готовности каждого результата контрольной суммы, в итоге мы получим последовательность чисел которая будет полностью лишена какого-либо смысла. Вместо результата первого потока мы сможем увидеть результат любого другого потока, только потому, что он вывел контрольную сумму на долю секунды раньше первого. А вторым примером общего доступа к ресурсам может стать индекс массива. Если потоки будут иметь общий доступ к индексам при умножении или любой другой операции над матрицами, мы получим неконтролируемый инкремент или декремент, что приведет к неправильной работе алгоритма, и соответственно к неправильному результату.

Примеры применения технологии инкрементного программирования при создании параллельного ПО

Рассмотрим на примере умножения двух матриц $C=A \times B$ как можно применить в конкретной программе, написанной на языке программирования C++, средства OpenMP с целью распараллеливания выполнения отдельных её частей.

Как известно, умножение двух матриц осуществляется по формуле:

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

Программа на языке программирования C++ будет иметь вид:

```

        c[i][j]=s;
    }
}

```

Из формулы умножения матриц мы видим, что вычисления промежуточных произведений являются независимыми, поскольку каждое из них записывает (и читает) свой элемент c_{ij} . Таким образом, их можно выполнять параллельно. Для этого в исходный код программы

вставляем соответствующие директивы OpenMP для параллельной реализации промежуточных произведений. Ниже приводится фрагмент текста программы на C++ с некоторыми вставками из OpenMP.

```

#include <omp.h> ...
...
void m_mult_m_OpenMP (double *a, double *b, double *c, int n)
    { int i, j, k;
  // запретить изменять количество потоков во время исполнения
//программы
  #pragma omp_set_dynamic(0);
  // установить количество потоков равным 2
  #pragma omp_set_num_threads(2);
    for(i = 0; i < n; i++)
        for(j = 0; j < n; j++)
            c[i][j]= 0;
  // описываем массивы, которые содержат элементы матриц, //глобально
видимыми, а индексы - локально видимыми
  #pragma omp parallel shared(a, b, c) private(i,j,k)
    #pragma omp for
    for (i = 0; i < n; i++){
        for (j = 0; j < n; j++){
            double s =0;
            for (k = 0; k < n; k++)
                s += a[i][k] * b[k][j];
            c[i][j]=s;
        }
    }
}

```

Результаты работы программы без распараллеливания и с распараллеливанием вычислений, т.е. с применением OpenMP, приведено в табл. 1. Программа выполнялась на 2-х ядерном процессоре. Конфигурация системы: DualCoe E6300 с 3.49ГБ ОЗУ, Windows XP.

Таблица 1. Время выполнения программы на 2-х ядерном процессоре

Размер матриц	Программа с OpenMP, мс	Программа без OpenMP, мс
200x200	32	63
500x500	890	1764
1000x1000	8766	15703
2000x2000	71469	124469
3000x3000	264328	425687

Из табл. 1 видно, что время решения задачи, использующей параллельные вычисления, в среднем на 50% меньше чем без.

В табл. 2 приведено время выполнения программы на одно и двух ядерном процессоре для матриц размером 3000x3000.

Таблица 2. Сравнение времени выполнения программы на одно и двух ядерных процессорах

Компьютер	Программа с OpenMP, мс	Программа без OpenMP, мс
Pentium4 3,0GHz +HT 2 ГБ ОЗУ	748156	730857
E6300 2x2,8GHz 3,49ГБ ОЗУ	320428	463157

Примечание: Pentium4- одноядерный процессор с реализацией технологии Hyper-threading (HT).

Из табл.2 видно, что время решения задачи на Pentium4 с применением OpenMP отличается незначительно от времени решения без ее применения, поскольку здесь не используется технология HT. Для применения этой технологии в полной мере необходимо использовать специальные API, описание которых можно найти на сайтах компаний Intel и Microsoft.

Для повышения эффективности программы умножения двух матриц, можно учитывать наличие кэш-памяти компьютера (кэш первого уровня). Когда при выполнении программы необходимо оперировать данными и они не находятся в кэш-памяти (кэш-промах), то процессор должен обращаться к оперативной памяти (ОП), что увеличивает временные затраты на выполнение задачи.

При проведении эксперимента было выявлено, что с ростом размера матриц теряется эффективность использования кэш-памяти, поскольку длина строк и столбцов матриц превосходят размеры кэш-памяти. Однако, если построить блочный алгоритм умножения матриц, то можно получить кратчайшее время выполнения этой операции: матрицы, которые умножаются, могут быть представлены как составленные из блоков, и результирующая мат-

рица будет получена путем выполнения умножения над блоками. Если размер блока соизмерим с размером кэш-памяти, то время вычисления произведения матриц уменьшится.

Рассмотрим еще один пример применения метода инкрементного программирования при создании ПО с использованием CUDA на графических процессорах для операции умножения двух матриц $C=A \times B$.

Как и в случае программирования с помощью технологии OpenMP, сначала необходимо выделить в алгоритме математические операции над массивами, которые могут быть распараллелены на графических процессорах, а затем шаг за шагом добавить новые команды и функции CUDA, описывающие параллельные конструкции.

В программах, написанных на языке Си с использованием технологии CUDA, помимо основной памяти CPU используется глобальная память (global), разделяемая память (shared) и другие.

Глобальная память – это память, которая выделяется на DRAM GPU. Для использования этой памяти программа должна с помощью функций CUDA SDK выделять (захватывать) память, выполнять двухстороннее копирование между CPU и GPU.

CUDA-функции выделения глобальной памяти под массивы имеют вид:

```
cudaMalloc ( (void**)&aDev, numBytes );
cudaMalloc ( (void**)&bDev, numBytes );
cudaMalloc ( (void**)&cDev, numBytes );
```

CUDA-функции копирования массивов данных из CPU на GPU имеют вид:

```
cudaMemcpy( aDev, a, numBytes, cudaMemcpyHostToDevice );
cudaMemcpy( bDev, b, numBytes, cudaMemcpyHostToDevice );
```

CUDA- функции копирования массивов данных из GPU на CPU имеют вид:

```
cudaMemcpy( c, cDev, numBytes, cudaMemcpyDeviceToHost );
```

Пример программы на языке Си с приведенными выше функциями CUDA которая функционирует на CPU, имеет вид:

```
...
#include <cuda_runtime.h>
#define BLOCK_SIZE 16
#define N 1024 //порядок матриц

int main ( int argc, char * argv [] )
{
    int numBytes = N * N * sizeof ( float );
    // выделение памяти под массивы на CPU
    float * a = new float [N*N];
    float * b = new float [N*N];
```

```

float * c = new float [N*N];
// инициализация массивов
for ( int i = 0; i < N; i++ )
    for ( int j = 0; j < N; j++ )
    {
        a [i] = 2.0f;
        b [i] = 1.0f;
    }
// выделение памяти под массивы на GPU
float * adev = NULL;
float * bdev = NULL;
float * cdev = NULL;
cudaMalloc ( (void**)&adev, numBytes );
cudaMalloc ( (void**)&bdev, numBytes );
cudaMalloc ( (void**)&cdev, numBytes );
// конфигурация сетки
dim3 threads ( BLOCK_SIZE, BLOCK_SIZE );
dim3 blocks ( N / threads.x, N / threads.y);
//старт cudaEvent функций, с помощью которых отслеживается
//завершение работы программ на GPU и определяется время //выполнения за-
дачи на GPU
cudaEvent_t start, stop;
float gpuTime = 0.0f;
cudaEventCreate ( &start );
cudaEventCreate ( &stop );
cudaEventRecord ( start, 0 );
// копирование массивов из CPU на GPU
cudaMemcpy ( adev, a, numBytes, cudaMemcpyHostToDevice );
cudaMemcpy ( bdev, b, numBytes, cudaMemcpyHostToDevice );
// Операция умножения подматриц выполняется
// в global-функции matrixMult
matrixMult <<<blocks, threads>>> ( adev, bdev, cdev, N );
// копирование массивов из CPU на GPU
cudaMemcpy ( c, cdev, numBytes, cudaMemcpyDeviceToHost );
//прекращение работы cudaEvent функций,
cudaEventRecord ( stop, 0 );
cudaEventSynchronize ( stop );
cudaEventElapsedTime ( &gpuTime, start, stop );
// печать времени выполнения задачи на GPU
printf("time spent executing by the GPU: %.2f milliseconds\n", gpu-
Time );
// освобождение ресурсов
cudaEventDestroy ( start );
cudaEventDestroy ( stop );
cudaFree ( adev );
cudaFree ( bdev );
cudaFree ( cdev );
delete a;
delete b;
delete c;
return 0;
}

```

Пример global-функции matrixMul, которая реализует массивно-параллельные операции умножения матриц с использованием global-памяти имеет следующий вид :

```

global__ void matrixMult(float *a, float *b, float *c)
{
    int bx = blockIdx.x; //индексы блока
    int by = blockIdx.y;
    int tx= threadIdx.x; //индексы нити внутри блока
    int ty= threadIdx.y;
    float sum =0.0f; //накопление результата

    // смещение для a[i][0]
    int ia = n * blockSize * by + n * ty;
    // смещение для b[0][j]
    int ib = blockSize * bx + tx;
    // умножаем и вычисляем сумму
    for (int k = 0; k < n; k++)
        sum += a[ia + k] * b[ib + k * n];
    // сохраняем смещение вычисленного элемента
    // в глобальной памяти
    int ic = n * blockSize * by + blockSize * bx;
    c[ic + n * ty + tx] = sum;
}
}

```

Однако global-память не имеет высокого быстродействия, чаще всего она используется для сохранения больших массивов данных, доступ к которым осуществляется как можно реже.

Для нахождения одного элемента результирующей матрицы нам необходимо выполнить чтение $2 \times N$ значений из global-памяти и исполнить $2 \times N$ арифметических операций, что является очень затратным по времени. Для уменьшения времени выполнения можно внести изменения в исходный код global-функции, применяя shared-память.

Shared-память представлена в виде блоков, которые находятся непосредственно в потоковом мультипроцессоре. Каждому блоку выделяется 16 Кбайт shared-памяти, доступ к которой и операции над массивами имеют очень высокое быстродействие. В этом случае результирующая подматрица C^* является суммой произведений подматриц [6]:

$$C^* = A_1^* \times B_1^* + A_2^* \times B_2^* + \dots + A_{N/16}^* \times B_{N/16}^*$$

В таком случае global-функции matrixMul, которая реализует этот алгоритм массивно-параллельных операций умножения матриц на shared-памяти имеет вид :

```

__global__ void matrixMult( float* a, float* b, float* c, int n)
{
    // Номер блока
    int bx = blockIdx.x;
    int by = blockIdx.y;
    // Номер нити
    int tx = threadIdx.x;
    int ty = threadIdx.y;
    // Индекс начала первой подматрицы A, которая обрабатывается блоком
    int aBegin = n * BLOCK_SIZE * by;
    // Индекс конца первой подматрицы A, которая обрабатывается блоком

    int aEnd = aBegin + n - 1;
    // Шаг перебора подматриц A
    int aStep = BLOCK_SIZE;
    // Индекс начала первой подматрицы B, которая обрабатывается блоком
    int bBegin = BLOCK_SIZE * bx;
    // Шаг перебора подматриц B
    int bStep = BLOCK_SIZE * n;
    float sum=0.0f; //Элемент, который вычисляется
    // Цикл по 16x16 подматрицам A и B

```

```

for (int ia = aBegin, ib = bBegin; ia <= aEnd; ia += aStep, ib += bStep)
{
    // Очередная подматрица A в shared-памяти
    __shared__ float as[BLOCK_SIZE][BLOCK_SIZE];
    // Очередная подматрица B в shared-памяти
    __shared__ float bs[BLOCK_SIZE][BLOCK_SIZE];
    // Загрузить по одному элементу из A и B в shared-память
    as[ty][tx] = a[ia + n * ty + tx];
    bs[ty][tx] = b[ib + n * ty + tx];
    // Дождаться, когда обе подматрицы будут подностью загружены
    __syncthreads();
    // Вычисляем элемент произведения загруженных подматриц
    for (int k = 0; k < BLOCK_SIZE; k++)
        sum += as[ty][k] * bs[k][tx];
    // Дождаться, когда все нити блока закончат вычисления
    __syncthreads();
}
// Записать результат в массив, который содержит матрицу c
int ic = n * BLOCK_SIZE * by + BLOCK_SIZE * bx;
c[ic + n * ty + tx] = sum;
}

```

Исходя из выше приведенного исходного кода вычисления произведения матриц, для определения одного элемента результирующей матрицы необходимо выполнить $2 \times N/16$ чтений из глобальной памяти и исполнить $2 \times N$ арифметических операций. Используя на компьютере графические процессоры, можно решить задачу в десятки раз быстрее. Можно сделать вывод, что использование shared-памяти существенно уменьшает суммарное время обращения к global-памяти.

Вывод

Время решения одной задачи на многоядерных ПК значительно сокращается при распараллеливании ПО с помощью технологии OpenMP или других средств параллельного программирования. Параллельное программирование – это технология будущего. Применение технологии инкрементного программирования упрощает работу пользователя как при

написании параллельных программ на OpenMP или CUDA, так и при их отладке. Кроме того, это дает возможность лучше понять суть параллельного программирования и подготовиться к работе на параллельных компьютерах кластерного типа.

Литература

1. <http://www.mpiforum.org/>
2. <http://www.epm.ornl.gov/pvm/>
3. <http://ru.wikipedia.org/wiki/openmp>
4. <http://www.openmp.org/>
5. «Основы работы с технологией CUDA» А.В. Боресков, А.А. Харламов, –М.: ДМК Пресс, 2010 год.
6. <http://developer.amd.com/gpu/ATIStreamSDK/Pages/default.aspx>
7. <http://ati.amd.com/products/streamprocessor/sp ecs.html>
8. Антонов А.С. Параллельное программирование с использованием технологии OpenMP, издательство МГУ, 2009