

УДК 681.03

А.В. Чистяков, И.С. Ислямова

О НЕКОТОРЫХ ОСОБЕННОСТЯХ РАЗРАБОТКИ ПРОГРАММНОГО ОБЕСПЕЧЕНИЯ ДЛЯ СОВРЕМЕННЫХ КОМПЬЮТЕРОВ

**Национальный
авиационный университет**

**Кафедра инженерии
программного
обеспечения**

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

Рассмотрены некоторые вопросы разработки прикладного программного обеспечения для современных компьютеров. Предложена технология планирования вычислений при создании программного обеспечения для решения вычислительных задач на компьютерах гибридной архитектуры. Приведены некоторые результаты компьютерных экспериментов.

Розглянуто труднощі розробки прикладного програмного забезпечення для сучасних комп'ютерів. Запропоновано технологію планування обчислень при створенні програмного забезпечення для розв'язування обчислювальних задач на комп'ютерах гібридної архітектури. Наведено деякі результати комп'ютерних експериментів.

An overview of some aspects the development of application software for modern computers. The technology of planning computer solutions by creating software to solve computational problems on computers of the hybrid architectures. Some results of computer experiments.

Ключевые слова: параллельное программирование, планирование вычислений, гибридные архитектуры, NVIDIA CUDA, GPGPU, OpenMP, MPI, PVM.

Введение

Современная наука находится в стадии динамичного развития. В каждой отрасли народного хозяйства возникают задачи, требующие обработки огромных объемов информации с высоким быстродействием и точностью. Объемы расчетов, необходимых для решения задач в различных сферах науки и промышленности постоянно растут, поэтому для их решения приходится строить все более мощные вычислительные системы.

Долгое время повышение производительности традиционных одноядерных процессоров в основном осуществлялось за счет последовательного увеличения тактовой частоты (около 80% производительности процессора определя-

ла именно тактовая частота) с одновременным увеличением количества транзисторов на одном кристалле. Однако дальнейшее повышение тактовой частоты (при тактовой частоте более 3,8 ГГц чипы просто перегреваются!), натывается на ряд фундаментальных физических барьеров. Кроме того, преимущества более высокой тактовой частоты сильно уменьшаются из-за задержек при обращении к памяти, так как время доступа к памяти не соответствует возрастающим тактовым частотам.

Типы современных компьютеров

Основным способом увеличения мощности современных компьютеров стало объединение все большего числа вычислительных элементов

в огромные вычислительные комплексы. Однако этот процесс в последние годы несколько замедлился из-за достижения пиковых возможностей технологии их производства. Недостатком такого увеличения вычислительных мощностей является то, что компьютеры получают все большие размеры, и как следствие становятся более сложными и гораздо более дорогими в эксплуатации, уже сегодня для них необходимо строить отдельные строения с развитой инфраструктурой, а потребление электроэнергии уже сопоставимо с потреблением небольшого городка.

Параллельно с совершенствованием центрального процессора (CPU) в последнее время происходит очень активное развитие графических процессоров (GPU) [1].

Ноябрь 2008 – фирма Intel представила линейку 4-ядерных центральных процессоров Intel Core i7, в основу которых положена микроархитектура нового поколения Nehalem. Процессоры работают на тактовой частоте 2,6 – 3,2 ГГц.

Декабрь 2008 – началась поставка 4-ядерного центрального процессора AMD Phenom II 940 (кодовое название – Deneb). Работает на частоте 3 ГГц.

Май 2009 – компания AMD представила версию графического процессора ATI Radeon HD 4890 с тактовой частотой ядра, увеличенной с 850 МГц до 1 ГГц. Это первый графический процессор, работающий на частоте 1 ГГц. Вычислительная мощность чипа, благодаря увеличению частоты, выросла с 1,36 до 1,6 терафлоп. Процессор содержит 800 (!) вычислительных ядер, поддерживает видеопамять GDDR5, DirectX 10.1, ATI CrossFireX и все другие технологии, присущие современным моделям видеокарт

Основными отличиями GPU по сравнению с CPU являются:

- архитектура максимально направлена на увеличение скорости расчета текстур и сложных графических объектов;
- пиковая мощность типичного GPU намного выше, чем у CPU;
- благодаря специализированной конвейерной архитектуре GPU намного эффективнее в обработке графической информации, чем CPU.

Сегодня наиболее известные типы компьютеров, которые используются для построения больших проектов, такие.

1). **Мэйнфреймы** или большие компьютеры – это компьютеры, созданные для обработки больших объемов информации. Наиболее круп-

ный их производитель – фирма IBM. Отличаются исключительной надежностью, высоким быстродействием, очень большой пропускной способностью каналов ввода-вывода. К ним могут присоединяться тысячи терминалов (дисплеев с клавиатурой) или персональных компьютеров для работы пользователей. Большинство крупных корпораций, банков, зарубежных правительственных учреждений обрабатывают свои данные именно на таких компьютерах. Хотя они могут стоить миллионы долларов, спрос на них не падает, так как централизованное хранение и обработка данных обходятся дешевле, чем обслуживание с помощью распределенных систем обработки данных, состоящие из сотен и тысяч персональных компьютеров.

2). **Суперкомпьютеры** – вычислительная машина, которая значительно превосходит по своим техническим параметрам большинство существующих компьютеров. Как правило, современные суперкомпьютеры представляют собой большое количество высокопроизводительных серверных компьютеров, соединенных друг с другом локальной высокоскоростной магистралью для достижения максимальной производительности в рамках распараллеливания вычислительной задачи. Эти компьютеры предназначены для решения задач, требующих огромных объемов вычислений. Основные потребители супер-ЭВМ – ученые, авиаторы, судостроители, военные, метеорологи, геологи, и многие другие. Супер-ЭВМ стоят десятки миллионов долларов (если не дороже), их производят лишь несколько крупных фирм, например, Cray Research (сейчас это подразделение фирмы Silicon Graphics), Hitachi и др. Такие компьютеры целесообразно использовать в глобальных системах, грид-среде и других компьютерных сетях.

3). **Рабочие станции** (в том числе с параллельной архитектурой) – персональные компьютеры, дополненные мощностями (комплекс технических средств, программное обеспечение, устройства), которые позволяют работать с большими объемами данных. Они имеют производительность, как у самых мощных персональных компьютеров или даже в несколько раз больше. Это компьютерная система, которая зачастую работает в составе компьютерной сети и настроена на выполнение задач инженеров, экономистов, программистов и других специалистов.

4). **Вычислительный кластер** – это совокупность компьютеров (вычислительных узлов с распределенной памятью), объединенных

единой коммуникационной сетью, используемых для решения одной или нескольких проектов. Современные вычислительные кластеры в мире могут насчитывать несколько тысяч вычислительных узлов. На сегодняшний день широко используется для решения научно-технических задач, в том числе в сетях.

5). **Гибридные системы** [2] – объединение нескольких вычислительных узлов, каждый из которых снабжен, как минимум, одним графическим устройством. Специалисты в разных развитых странах мира отмечают, что именно за таким компьютером они видят будущее компьютеризации.

Системы программирования для решения математических задач

Появление мощных компьютеров новой архитектуры привело к созданию новых средств программирования. Наиболее популярными компьютерами для высокопроизводительных вычислений являются многоядерные (компьютеры MIMD-архитектуры): каждый узел у них есть система с распределенной памятью. Для создания программного обеспечения на такие компьютеры появились специальные программные средства PVM (Parallel Virtual Machine) [3], MPI (Message Passing Interface) [4], [5] и другие.

PVM общедоступный программный пакет, позволяющий объединять разнородный набор компьютеров в общий вычислительный ресурс («виртуальную параллельную машину») и предоставляющий возможности управления процессами с помощью механизма передачи сообщений. PVM – это продукт продвижения гетерогенных сетевых исследовательских проектов, распространяемый авторами и институтами, в которых они работают. Поддерживается языками программирования Фортран, Си, Си++.

Чаще всего для написания программ для компьютеров с распределенной памятью используется система MPI – стандартизированный интерфейс для построения программ по модели обмена сообщениями. Существуют бесплатные и коммерческие реализации для различных суперкомпьютерных платформ.

Преимущества:

- возможность использования в языках Фортран, Си, Си++, возможность совмещения обменов сообщениями и вычислений, несколько режимов передачи сообщений;

- широкий набор коллективных операций;

- широкий набор редуцированных операций;

- удобные средства именованного адресатов сообщений;

- возможность задания типа передаваемой информации.

Недостатки:

- слишком сложный для прикладного программиста (программирование на низком уровне, необходимо знать: как поделить между процессорами работу программ, как поделить между процессорами данные, как организовать доступ к удаленным данным или обмен данными, которые в процессе вычислений нужны одному процессору, в то время как хранятся на другом).

Для программирования многопоточных приложений на многопроцессорных системах с общей памятью (SMP-системах) используется система OpenMP (Open Multi-Processing) [6]. OpenMP – это система программирования, которая расширяет последовательную программу (Фортран, Си, Си++) набором директив распараллеливания. Предполагается, что потоки выполняются параллельно на машине с несколькими процессорами (количество процессоров не обязательно должно быть больше или равно количеству потоков). Задачи, выполняемые потоками параллельно, также как и данные, требуемые для выполнения этих задач, описываются с помощью специальных директив препроцессора соответствующего языка – *прагм*. Программа начинает свое выполнение как один процесс (главная нить), пока не встретится параллельная область программы (ограничивается директивами *PARALLEL* и *END PARALLEL*). При входе в параллельную область главная нить порождает некоторое число подчиненных ей нитей, образуя группу нитей. Все операторы программы, находящиеся в параллельной конструкции, выполняются всеми нитями группы параллельно, пока не произойдет выход из параллельной области или не встретится одна из конструкций распределения работ. При выходе из параллельной конструкции все порожденные ранее нитки сливаются с главной, которая продолжает выполняться дальше последовательно. “Инкрементальное распараллеливание” программы, которое используется в OpenMP идеально подходит для разработчиков, желающих быстро распараллелить свои вычислительные программы с большими параллельными циклами. Разработчик не создает новую параллельную программу, а просто последовательно добавляет в текст последовательной программы соответствующие OpenMP-директивы.

В последние годы стали активно использоваться графические процессоры для неграфических вычислений (моделирование физических процессов, обработка сигналов, вычислительная математика / геометрия, операции с базами данных, вычислительная биология, вычислительная экономика, компьютерная визуализация и т. д.). Однако при использовании графических процессоров для неграфических задач были значительные трудности, связанные с тем, что не было никакого стандартного интерфейса для программирования. Разработчики использовали OpenGL или Direct3D, но это было очень неудобно, потому что приходилось знать и учитывать много технических особенностей вычислительной среды графических ускорителей.

Устройство GPU отличается от того, как устроен CPU. Прежде всего, отличается работа с памятью. Так, не все центральные процессоры имеют встроенные контроллеры памяти, а у всех GPU обычно есть по несколько контроллеров, вплоть до восьми 64-битных каналов в чипе NVIDIA GT200. Кроме того, на видеокартах применяется более быстрая память, и в результате видеочипам доступна в разы большая пропускная способность памяти, что весьма важно для параллельных расчетов, оперирующих с огромными потоками данных.

В универсальных процессорах большие количества транзисторов и площадь чипа идут на буферы команд, аппаратное предсказание ветвления и огромные объемы кэша памяти. Все эти аппаратные блоки нужны для ускорения исполнения немногочисленных потоков команд. Видеочипы тратят транзисторы на массивы исполнительных блоков, управляющие потоками блоков, разделяемую память небольшого объема и контроллеры памяти на несколько каналов. Вышеперечисленное не ускоряет выполнение отдельных потоков, оно позволяет чипу обрабатывать нескольких тысяч потоков, одновременно исполняющихся чипом и требующих высокой пропускной способности памяти.

Кроме того, есть существенные отличия в кэшировании. Универсальные центральные процессоры используют кэш-память для увеличения производительности за счет снижения задержек доступа к памяти, а GPU используют кэш или общую память для увеличения полосы пропускания. CPU снижают задержки доступа к памяти при помощи кэш-памяти большого размера, а также предсказания ветвлений кода. Эти аппаратные части занимают большую часть площади чипа и потребляют много энер-

гии. Видеочипы обходят проблему задержек доступа к памяти при помощи одновременного исполнения тысяч потоков – в то время, когда один из потоков ожидает данных из памяти, видеочип может выполнять вычисления другого потока без ожидания и задержек.

Есть множество различий и в поддержке многопоточности. CPU исполняет 1-2 потока вычислений на одно процессорное ядро, а видеочипы могут поддерживать до 1024 потоков на каждый мультипроцессор, которых в чипе несколько штук. И если переключение с одного потока на другой для CPU стоит сотни тактов, то GPU переключает несколько потоков за один такт.

Кроме того, центральные процессоры используют SIMD (одна инструкция выполняется над многочисленными данными) блоки для векторных вычислений, а видеочипы применяют SIMT (одна инструкция и несколько потоков) для скалярной обработки потоков. SIMT не требует, чтобы разработчик преобразовывал данные в векторы, и допускает произвольные ветвления в потоках.

Вкратце можно сказать, что в отличие от современных универсальных CPU, видеочипы предназначены для параллельных вычислений с большим количеством арифметических операций. И значительно большее число транзисторов GPU работает по прямому назначению — обработке массивов данных, а не управляет исполнением (flow control) немногочисленных последовательных вычислительных частей программы.

В 2005–2006 г. корпорация NVIDIA (один из крупнейших производителей графических, медиа- и коммуникационных процессоров) начала разработку единого и удобного стандарта, и в 2007 г. представила первоначальную версию технологии CUDA для программирования на GPU. Разработчики ставили цель, чтобы GPU использовались не только для создания изображений в 3D приложениях, но и применялись в других параллельных расчетах.

Технология NVIDIA CUDA [7] – архитектура CUDA основана на концепции одна команда на множество данных (Single Instruction Multiple Data, SIMD) и понятии мультипроцессора. Концепция SIMD подразумевает, что одна инструкция позволяет одновременно обработать множество данных. Мультипроцессор – это многоядерный SIMD процессор, позволяющий в каждый определенный момент времени выполнять на всех ядрах только одну инструкцию. Каждое ядро мультипроцессора скалярное, т.е. оно не поддерживает векторные

операції в чистому вигляді. В CUDA є поняття пристрою (*device*), що означає програма на відеокарті, підтримувач драйвер CUDA, а також – хост (*host*), т.е. програма в звичайній оперативній пам'яті комп'ютера, використовуюча CPU і виконуюча управляючі функції по роботі з пристроєм.

CUDA передбачає спеціальний підхід к розробці, не зовсім такої, як прийнято в програмах для CPU. Потрібно пам'ятати про різні типи пам'яті, про те, що локальна і глобальна пам'ять не кешується і затримки при доступі до неї значно вище, ніж у регістрової пам'яті, так як вона фізично знаходиться в окремих мікросхемах. Особливістю архітектури CUDA є блочно-сіткова організація, незвичайна для багатопотокових програм (рис. 1). При цьому драйвер CUDA самостійно розподіляє ресурси пристрою між потоками. На рис. 1. ядро позначено як *Kernel*.

CUDA програми програмується на C++ з використанням спеціального SDK і представляє собою комбінацію коду, який виконується на CPU і ядер (*kernels*), які виконуються на GPU. Кожна задача ділиться на незалежні частини – підзадачі (подібно тому, як робиться при допомозі OpenMP), які вирішуються з допомогою набору взаємодіючих між собою ниток. CUDA використовує велику кількість окремих ниток для обчислень, часто кожному обчислюваному елементу відповідає одна нить. Всі нитки групуються в ієрархію – *grid*, *block*, *thread*. Розподіл обчислень і даних повинен бути узгодженим і виконуватися з урахування чергування послідовних і паралельних частин програми. Верхній рівень – *grid* – відповідає ядру і об'єднує всі нитки, які виконують дану ядро. *Grid* представляє собою одномерний або двохмерний масив блоків (*block*). Кожен блок є одно / двох / трьохмерний масив ниток (*threads*). При цьому кожен блок – повністю незалежний набір взаємодіючих між собою ниток, нитки з різних блоків не можуть між собою взаємодіювати.

Виконання розрахунків на GPU показує чудові результати в алгоритмах, використовують паралельну обробку даних. Тобто, коли одну і ту ж послідовність математичних операцій застосовують до великого обсягу даних. При цьому кращі результати досягаються, якщо відношення числа арифметичних інструкцій до числа звернень до пам'яті достатньо велике. Це пред'яв-

ляє менші вимоги до управління виконанням, а висока щільність математики і великий обсяг даних скасовує необхідність в великих кешах, як на CPU. При великому кількості GPU система CUDA дозволяє прискорити рішення задачі в сотні разів.

На сьогоднішній день вже створено багато бібліотек програм для математичних розрахунків, особливо по чисельній математиці. В основному в них використовується CUDA версії 2.1. Одним з найбільш великих недоліків в цій системі програмування є те, що зв'язок між окремими GPU підтримується тільки через CPU, що суттєво уповільнює роботу програми і ускладнює процес програмування.

Однак в лютому 2011 року компанія NVIDIA анонсувала нову версію набору інструментів NVIDIA CUDA 4.0 для розробки паралельних програм з допомогою графічних процесорів. Новий набір інструментів підтримує технологію NVIDIA GPU Direct 2.0, яка забезпечує рівноправну зв'язок між GPU в межах одного сервера або робочої станції, що спрощує і прискорює мультипроцесорне програмування і роботу програм. Уніфікована віртуальна адресація (UVA) організує єдине адресне простір для основної системної пам'яті і пам'яті GPU, що робить паралельне програмування ще швидше і простіше. Набір також включає бібліотеки примітивів Thrust C++ – структур даних і алгоритмів C++ з відкритим кодом для паралельних розрахунків.

Для ефективного рішення задач на комп'ютерах гібридної архітектури, які можуть нарахувати декілька сотень багатоядерних CPU і GPU, створюються спеціальні алгоритми і програми по гібридній технології [8]. Гібридна технологія – це два рівні паралелізму – паралелізм задачі на підзадачі і паралелізм всередині підзадачі. В разі об'єднання декількох чисельних багатоядерних вузлів з декількома графічними пристроями задача розпаралелюється на окремі підзадачі шляхом розподілу між ядрами процесора елементів масивів і витків циклів з допомогою функцій бібліотеки OpenMP, які в свою чергу розпаралелюються на GPU. Обмін даними між окремими вузлами виконується з допомогою функцій бібліотеки MPI: організується взаємодія підзадач – обмін між ними граничними значеннями. Основний недолік підходу: програмісту потрібно знати і уміти

использовать сразу несколько разных моделей параллелизма и разные инструментальные

средства.

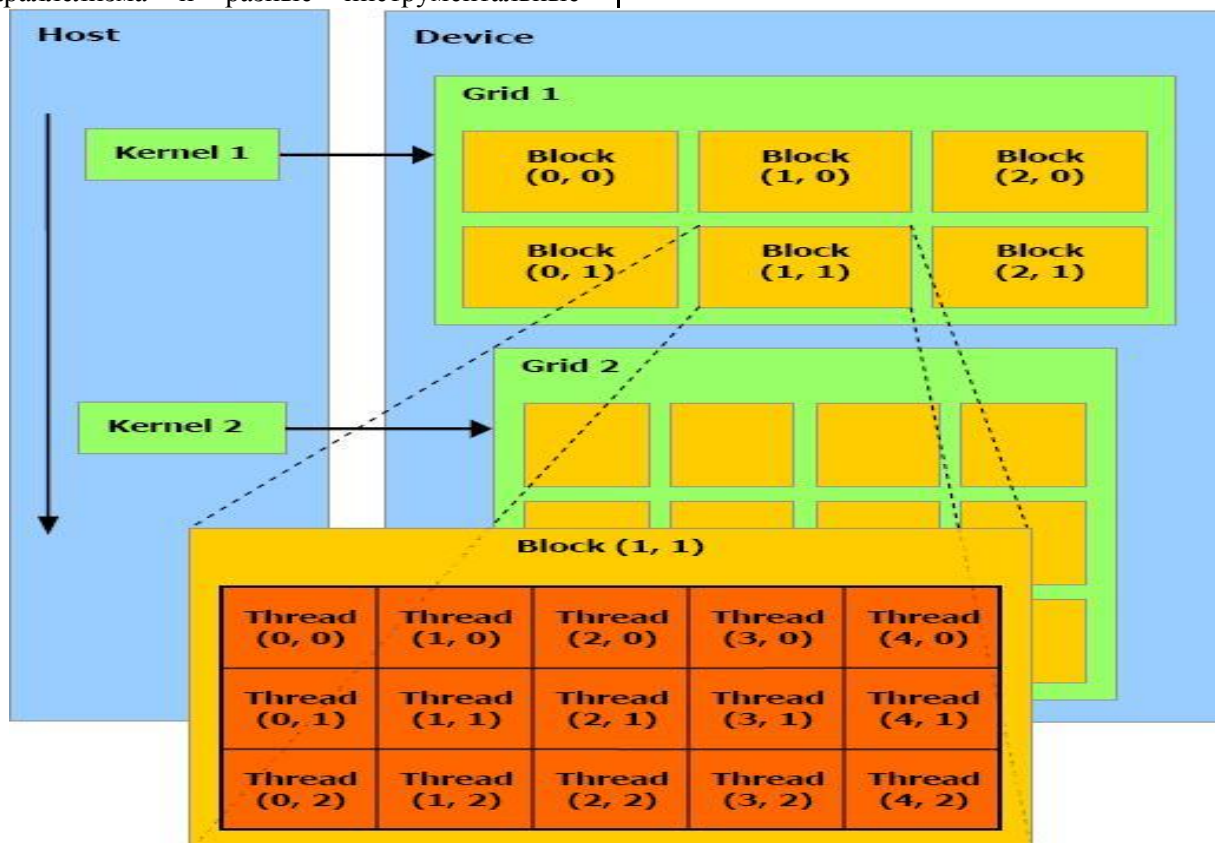


Рис. 1. Организация потоков при использовании технологии CUDA

Таким образом, создание алгоритмов и программ для решения задач с использованием той или иной архитектуры, которая может быть наиболее эффективной по времени выполнения задачи, является достаточно непростой работой. Пользователю необходимо не только обладать навыками программирования MPI, OpenMP, CUDA и т.д., но также хорошо знать архитектуру компьютера, вычислительные возможности CPU и GPU для оптимального выбора количества процессов, ядер, подключения графического ускорителя с целью скорейшего получения решения задачи с высокой точностью.

Планирование вычислений при разработке прикладного программного обеспечения

Вместе с бурным развитием вычислительной техники бурно развивается область компьютерной программной технологии (software engineering), которая требует периодического пересмотра подходов к созданию прикладного программного обеспечения.

Успешное использование CPU и GPU, высокая производительность компьютера зависит от успешной разработки вычислительного

программного обеспечения. В последнее время все чаще отмечается, что для облегчения работы пользователей на сложных вычислительных комплексах необходимо создавать библиотеки и пакеты прикладных программ с некоторым дружественным интерфейсом по детальной планировке вычислений на гибридных системах с правильным выбором ее вычислительных компонентов. Сначала необходимо разработать эффективные алгоритмы и программы для одного GPU и CPU, а также для нескольких этих ресурсов одновременно (гибридное задание), после чего определить ключевые параметры задачи и топологию компьютера, при которых решение задачи будет получено с самым высоким быстродействием и с оптимальным количеством необходимых вычислительных ресурсов.

Используя созданное программное наполнение и имея знания об условиях использования его компонент на различных архитектурах гибридной системы, можно разрабатывать прикладное программное обеспечение с автоматическим планированием вычислительных средств, что значительно повысит производительность работы как пользователя, так и ком-

пьютера. При цьому необхідно надавати користувачеві відповідні рекомендації щодо доцільності використання компонентів гібридної системи для рішення конкретної задачі.

Покажемо на простому прикладі множення двох матриць $C = AB$, як залежить час рішення задачі від її об'єму та архітектури комп'ютера. Як відомо, саме ця математична операція частіше за все виконується при рішенні прикладних задач і на її виконання витрачається багато комп'ютерного часу (n^3 - операцій, n - порядок матриць).

На рис. 2 схематично показано розпаралелювання задачі на гібридній системі, що має три CPU та три GPU.

Нагадаємо, що множення двох квадратних матриць виконується за формулою:

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

Для ефективного використання обчислювальних вузлів комп'ютера гібридної архітектури при виконанні цієї математичної операції можна організувати планування обчислень за такою схемою.

1). Якщо розміри матриць такі, що матриці повністю розміщуються в кеш-пам'яті комп'ютера, то задачу (1) можна розв'язувати на CPU (один процес) звичайним алгоритмом, що складається з трьох циклів.

2). Якщо розміри матриць такі, що матриці не поміщаються в кеш-пам'ять, то доцільно використовувати блокові алгоритми [9], при цьому необхідно визначити розміри бло-

ків матриць такі, щоб вони повністю розміщались в кеш-пам'яті.

3). Для кожного комп'ютера програму блокового множення матриць можна зробити більш ефективною, якщо виконати розкладку внутрішнього циклу з урахуванням конвейєра інструкцій комп'ютера.

4). Якщо CPU - багатоядерний процесор, то можна розпаралелити задачу на підзадачі між ядрами, використовуючи OpenMP (комунікаційні втрати не повинні перевищувати час виконання арифметичних операцій).

5). Якщо CPU має графічний прискорювач, то доцільно організувати обчислення за програмою з використанням CUDA: основна частина програми (визначення вихідних даних, управління GPU, обмін інформацією між CPU та GPU) буде виконуватися на CPU (*хост*), а паралельна частина програми (виконання циклів) буде виконуватися на GPU (*device*). При цьому час виконання задачі може зменшитися в десятки разів (для великих розмірів матриць і при великій кількості графічних прискорювачів - в сотні разів).

6). Якщо комп'ютер складається з кількох CPU та GPU, то необхідно передбачити програму з розпаралелюванням за двошаровою моделлю: наперед визначити оптимальну кількість процесорів для найшвидшого рішення задачі. Потім розпаралелити задачу на підзадачі за допомогою MPI, при цьому виконання циклів кожної підзадачі на GPU організувати за допомогою CUDA. Запустити задачу на виконання на встановлену кількість процесорів.

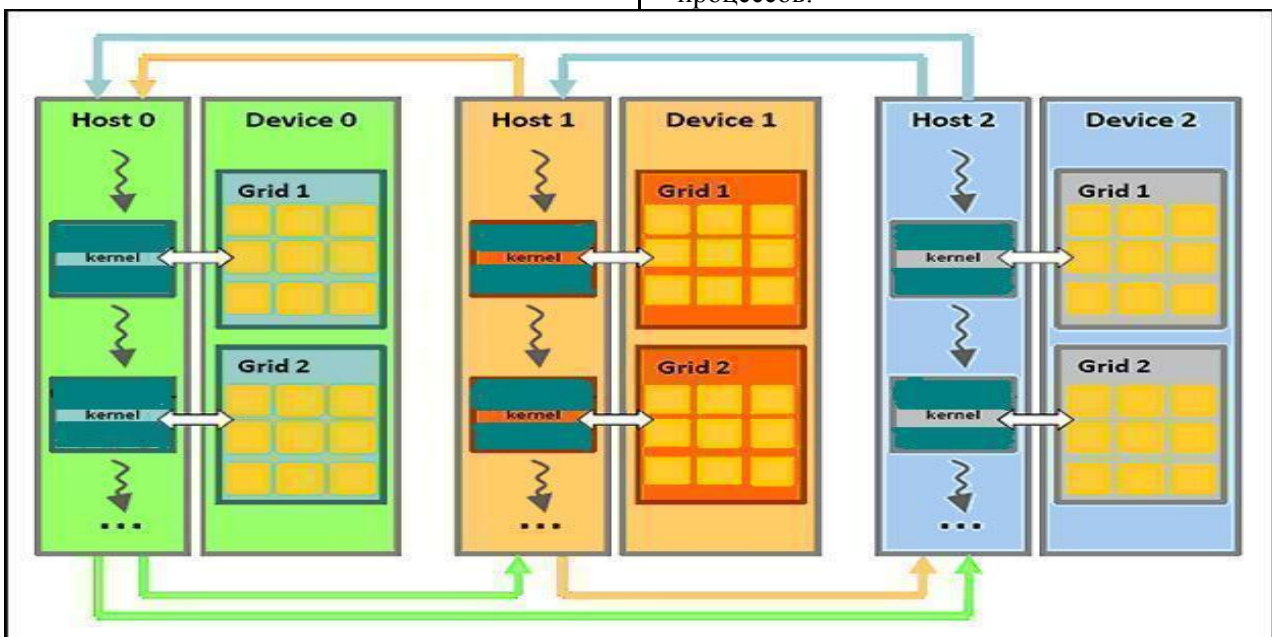


Рис. 2. Схема множення матриць на комп'ютері з трьома CPU і трьома GPU

В табл.1-3 демонструються результати перемноження матриць різними варіантами на гібридному комп'ютері, що складається з двох вичислювальних вузлів, з'єднаних комунікаційною мережею Infiniband. В вичислювальному вузлі використовується два чотириядерних процесора Intel Core i7 930. Тактова частота ядра становить 2.8 ГГц. На кожному вузлі є графічна відеокарта GeForce GTX480. Час виконання завдань наведено в сек.

В табл. 1 представлені результати рішення задачі на одному процесорі трьома вичислювальними схемами з перших трьох варіантів. С найбільшою швидкістю отримали рішення на одному процесорі блочним алгоритмом (розмір блоку 80) з розкладом внутрішнього циклу.

Із табл. 2 видно, що час рішення задачі (починаючи з порядку 1000), який розпаралелено на декілька ядер на процесорі з допомогою системи OpenMP, на двох ядрах декілько більше ніж звичайним способом на одному ядрі (комунікаційні витрати достатньо великі порівняно з часом виконання арифметики). Однак, починаючи з трьох ядер, задача для всіх представлених порядків вирішується значно швидше.

В табл. 3 показано як покращується час рішення задачі на чотирьох процесорах з ростом порядку матриць по паралельній програмі з використанням MPI+CUDA (задача розпаралелювалась на CPU і GPU) порівняно з часом рішення по паралельній програмі з використанням MPI (задача розпаралелювалась тільки на CPU).

Таблиця 1. Час множення матриць варіантами 1-3

Порядок матриць	Варіант 1	Варіант 2	Варіант 3
500	1	0.23	0.1
1000	13	5	3
2000	118	54	38
3000	452	128	61

Таблиця 2. Час множення матриць на різному кількості ядер комп'ютера

Порядок матриць	Звичайний алгоритм (1 ядро)	Алгоритм з використанням OpenMP		
		2 ядра	3 ядра	4 ядра
500	1	0.8	0.61	0.44
1000	13	18.9	7.4	4.01
2000	118	160	54	32
3000	452	612	255	109

Таблиця 3. Час множення матриць на 4 процесорах з використанням MPI і MPI+CUDA

Порядок матриць	Час рішення на CPU	Час рішення на CPU+GPU
1000	2.8	4.7
2000	25	5.14
3000	90	6.35

Висновки

В наше час, коли архітектура комп'ютерних комплексів дуже ускладнилась, ускладнюється також і програмне забезпечення, яке повинно враховувати як параметри задачі, так і умови ефективного рішення на комп'ютерній системі.

Сучасною тенденцією в світі інформаційних технологій є перехід до нової парадигми прикладного програмного забезпе-

чення – автоматичне планування рішення задачі з метою найменших витрат комп'ютерного часу при оптимальному використанні багатоядерних GPU і CPU.

Багато авторів відзначають, що для комп'ютерних систем сьогодні особливо актуально створення прикладного програмного забезпечення з автоматичним вибором вичислювальних ресурсів, алгоритмів і програм рішення задач.

Как следует из приведенных результатов экспериментов, использование вычислительных ресурсов гибридной системы дает значительные преимущества по эффективности вычисления (уменьшение времени решения задачи в десятки раз), если они учитывают как параметры задачи, так и вычислительные ресурсы. С ростом количества многоядерных процессоров и графических ускорителей время решения больших задач может уменьшиться в сотни раз.

Литература

1. Графические процессоры для высокопроизводительных вычислений: <http://gpu.parallel.ru/>

2. Свенч А.А., Файзуллин Р.Т., Хныкин И.Т. Гибридная суперкомпьютерная система. // Сборник трудов Международной научной конференции «Параллельные вычислительные технологии 2011», Московский государственный университет имени М.В. Ломоносова, Москва, 2011 г. (http://kvap.ru/pavt2011_cd/full.html).

3. Немнюгин С.А., Стесик О.Л. Параллельное программирование для многопроцессорных вычислительных систем. – СПб.: БХВ-Петербург, 2002. – 400 с.

4. К.Ю. Богачёв Основы параллельного программирования. – М.: БИНОМ. Лаборатория знаний, 2003. – 342 с.

5. Корнеев В.Д. Параллельное программирование в MPI. - Новосибирск: Изд-во СО РАН, 2000. – 213 с.

6. Чистяков А.В., Ислямова И.С. Метод и технологии параллельного программирования при решении прикладных задач. // Инженерия программного обеспечения. № 3 – 2010. – С.37 – 50.

7. Боресков А.В., Харламов А.А. Основы работы с технологией CUDA. – М.: ДМК Пресс, 2010. – 232 с.

8. Chorley M.J., Walker D.W. Performance analysis of a hybrid MPI/OpenMP application on multi-core clusters. Journal of Computational Science 1(3). – 2010. – P. 168–174.

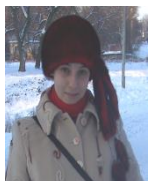
9. Fox G.C., Otto S.W. and Hey A.J.G. Matrix Algorithms on a Hypercube I: Matrix Multiplication Parallel Computing. – 1987. – 4. – P. 17–31.

Сведения об авторах



Чистяков Алексей Валериевич – студент 5 курса Национального авиационного университета. Научные интересы - параллельное программирование, прикладное и объектно-ориентированное программирование, инженерия программного обеспечения, компьютерные сети.

E-mail: bmwwmb@gala.net



Ислямова Инна Сергеевна – студентка 5 курса Национального авиационного университета. Научные интересы - параллельное программирование, прикладное и объектно-ориентированное программирование, объектно-ориентированные базы данных, инженерия программного обеспечения.

E-mail: Inna_Islyamova@ukr.net



Иванова Любовь Николаевна – к.т.н., доцент кафедры инженерии программного обеспечения Национального авиационного университета, научные интересы - параллельное программирование, инженерия программного обеспечения.

E-mail: lubov.ivanova@livenau.net