

УЧЕБНАЯ ПРОГРАММА

«Многопоточные вычисления на основе технологий CUDA и OpenCL»

Лекции: 26 часов

Практические занятия: 24 часов

Самостоятельная работа: 22 часа

Всего часов: 72

Лекторы: аспирант МФТИ А.М. Казеннов, О.В. Геллер.

Ассистент: студент МФТИ А.Е. Алексеенко.

1. Введение в курс

История развития вычислительных систем. Основная терминология курса. Типы параллелизма. Обоснование необходимости использования распределенных вычислительных систем. Критерии применимости параллельных вычислений. Примеры применения параллельных вычислений. Различные типы параллельных систем. Классическая и гибридная схема. Кластеры и суперкомпьютеры на гибридной схеме.

2. Архитектура CPU и GPU

Сравнение классической архитектуры Intel и AMD. Принципиальное отличие классической и GPU архитектуры. Необходимые шаги к единой архитектуре вычислительных устройств. Сравнительные характеристики чипов G 280, G 295, G 480, NVIDIA.

3. Аппаратная реализация единой архитектуры

Объединённая архитектура графических процессоров. Основные составные элементы аппаратной реализации GPU. Преимущества унифицированной архитектуры. Составные части аппаратной реализации: TPC, SM, SP. Буфер инструкций SM. Регистровый файл SM. Конвейеры исполнения команд. Ветвление внутри варпа.

4. Программная модель CUDA

Основные модификаторы языка C. Введение в особенности программирования под GPU. Понятия Thread, Warp, Block и Grid. Программный стек CUDA. Описание пользовательского интерфейса разработчика, основные компоненты. Команды работы с памятью. Пример вызова CUDA.

5. Программная модель OpenCL

Понятия Host и Device. Платформы OpenCL, контекст и очередь исполнения. Сборка и запуск ядер на устройствах.

6. Язык OpenCL C

Основные модификаторы языка C. Введение в особенности программирования под GPU. Понятия Work Item, Work Group и Warp. Описание пользовательского интерфейса разработчика, основные компоненты. Команды работы с памятью. Пример вызова OpenCL. Компиляция и запуск OpenCL программ.

7. Первые программы на CUDA и OpenCL

Необходимые приложения для написания программы. Установка программного обеспечения. Использование Putty и WinSCP для доступа к серверу. Написание программы нахождения числа Пи. Модификация для нахождения центра масс.

8. Модель памяти GPU

Глобальная, константная, текстурная, локальная, разделяемая и регистровая память. Особенности использования каждого типа памяти. Размещение различных данных в различной памяти. Сравнения производительности глобальной и текстурной памяти на задачах произвольного чтения. Характерные размеры каждой памяти на примере чипа G200. Когерентное общение с глобальной памятью.

9. Глобальная и разделяемая память GPU

Написание программы перемножения матриц с использованием глобальной памяти. Оптимизация написанной программы с использованием разделяемой памяти. Способы избегания конфликта банков в разделяемой памяти. Постановка практических заданий.

10. Оптимизация основных алгоритмов

Использование Scan, Reduce, Histogram, Bitonic sort.

11. Текстурная память GPU

Использование текстурной памяти. Способы размещения данных в текстурной памяти. Использование аппаратной интерполяции. Отличия модели исполнения, работы с текстурами, сборки и компиляции программ OpenCL от CUDA. Постановка практических заданий.

12. Итоговый проект

Постановка и разбор проектных заданий. Консультации по проектам. Прием заданий и проектов.

Примеры тем проектных работ

1. Клеточные автоматы типа «Жизнь».
2. Клеточные автоматы Кохомото-Ооно.
3. Решение простых двумерных сеточных итерационных задач.

Литература

1. Архитектура и программирование массивно параллельных процессоров: http://www.nvidia.ru/object/cuda_state_university_courses_ru.html
2. GPU Gems 1, 2, 3 edited Hubert Nquyen from NVIDIA.
3. Спецификация OpenCL <http://www.khronos.org/registry/cl/specs/opencl-1.1.pdf>
4. Introduction_to_OpenCL_Programming. <http://developer.amd.com>
5. Боресков А.В., Харламов А.В. Основы работы с технологией CUDA. – Изд-во: ДМК Пресс, 2010, 232 стр.