

Реализация Multi-GPU версии OpenFOAM

Александр Монаков, amonakov@ispras.ru

Институт Системного Программирования

6 декабря 2012 г.

Массивно-параллельные вычислительные устройства

- Ресурсы производительности на порядок выше по сравнению с процессорами:
 $\approx 1\text{--}3 \text{ TFLOP/s}$, $\approx 100\text{--}250 \text{ GB/s}$
- Акселераторы используют собственную память, доступную для CPU через PCI-Express
- Существенный прирост общей производительности и сокращение энергозатрат для подходящих задач
- Как правило, требуют использование специализированных структур данных и алгоритмов для достижения оптимальной производительности

- Распараллеливание в OpenFOAM достигается только за счёт разделения сетки между MPI-процессами
- Наиболее затратным является решение симметричных систем (50 – 80%)

Перенос метода сопряжённых градиентов на GPU:

- Работа с разреженной матрицей системы на GPU
- Предобуславливание
- Распараллеливание солвера

- Единственная операция, требующая оптимизации – умножение на плотный вектор
- Существенно зависит от выбора структуры данных
- Мы используем формат Sliced ELLPACK
 - Более высокая производительность по сравнению с библиотеками CUSPARSE, cusp
 - Требуют больше времени для преобразования формата

Как правило, портрет разреженной матрицы одинаков для всех решаемых систем

- Можно выполнить преобразование формата один раз
- Запомнить портрет и далее обновлять только ненулевые коэффициенты на GPU
- Важная оптимизация для сложных форматов

- Эффективное предобуславливание для GPU-акселераторов – активно исследуемая тема
- Преобуславливатели ILU-семейства не подходят для GPU:
 - Решение треугольных систем плохо параллелится
 - Трудности с векторизацией кода
- Удобно использовать предобуславливатели, основанные на построении приближённой обратной матрицы явно

Для реализации был выбран AINV, представляющий обратную матрицу в виде произведения треугольных:

$$W^T A Z = D + e$$

W , Z – верхнетреугольные с единицами на диагонали, $W = Z$ для симметричных
Построение параметризовано drop tolerance

- Применение на GPU – две операции умножения разреженной матрицы на вектор
- Оптимизация drop tolerance – открытый вопрос

Временные затраты на вычисление AINV на CPU соизмеримы с временем решения системы на GPU

Мы используем оптимизацию, основанную на предположении, что последовательно решаемые системы имеют близкие коэффициенты, и преобуславливатели от ранее решённых систем хорошо подходят для новых

- Можно несколько раз переиспользовать предобуславливатель
- Можно полностью скрывать затраты, вычисляя его асинхронно и параллельно с решением системы на GPU

- Хранить и применять предобуславливатель в одинарной точности (float)
- Переупорядочивание множителей в лево-верхнетреугольную форму

- 1 узел, несколько MPI-процессов + 1 GPU
 - Промежуточный этап между последовательной и параллельной реализациями
 - Полезный режим для исследователей с компьютерами с GPU
- Несколько узлов, $N_{CPU_s} > N_{GPU_s}$

- GPU используется мастер-процессом
- Без MPI-коммуникаций в солвере
- Решается объединённая матрица системы

Цель – эффективная реализация GPU-решателя для случая

$$N_{CPU_s} \geq N_{GPU_s}$$

... без одновременного использования GPU из разных процессов

- Количество MPI-процессов = количество CPU-ядер
- Каждый из доступных акселераторов используется лишь одним MPI-процессом

Преимущества схемы:

- Увеличение эффективности
- Хорошая совместимость

Использование CUDA+MPI в плагине:

- 1 GPU CUBLAS → «pinned»-память → MPI-редукция
- 2 GPU ядро ↔ оперативная память ↔ MPI Send/Recv

Простая схема, не требует поддержки CUDA в реализации MPI

17 узлов: 2x 6-core Xeon 5650, 3x Tesla M2070

51 GPU, 204 ядра CPU

Сетка с 50М ячеек

Только ресурсы CPU:

- ① 51 ядро: 2553 s
- ② 216 ядер: 1314 s

С использованием GPU:

- ① 51 ядро + 51 GPU: 633 s
- ② 102 ядра + 51 GPU: 525 s
- ③ 204 ядра + 51 GPU: 531 s