

## О ПРИМЕНЕНИИ GPU NVIDIA АРХИТЕКТУРЫ KEPLER

Кузьминский М.Б., Андреев А.М.

Институт органической химии им. Н.Д. Зелинского РАН, г. Москва

Кузьминский М.Б.  
Андреев А.М.  
Институт органической  
химии им. Н.Д.  
Зелинского РАН

GPU Nvidia архитектуры Kepler (K20, K20X, K40) являются мировыми лидерами производительности с плавающей запятой. В работе дан краткий анализ архитектурно-технических характеристик GPU Kepler на базе процессора GK110 и приведены некоторые данные об их производительности, в т.ч. полученные авторами.

Пиковая производительность этих GPU лежит в диапазоне 1.2-1.4 TFLOPS с двойной точностью (до 4.3 TFLOPS с одинарной точностью) в зависимости от модели, емкость главной памяти – от 5 до 12 Гбайт, пропускная способность – от 208 до 288 Гбайт/с. GPU подсоединяются к хосту через разъем PCI-E v.2 x16, за исключением K40, где применяется в 2 раза более быстрый канал PCI-E v.3 x16.

Пиковая производительность GPU на типичной для HPC-приложений двойной точности в разы превосходит пиковую производительность типичных для кластеров двухпроцессорных серверов при использовании в последних самых мощных серверных процессоров x86-64 Intel и AMD. Поэтому измерения производительности GPU Kepler представляет большую актуальность.

Такие измерения производительности для HPC можно отнести к одному из трех уровней: низшему, где производятся измерения отдельных технических характеристик GPU; среднему – где измерения относятся к универсальным математическим алгоритмам (например, умножение матриц или тесты Linpack); и к верхнему – на котором измеряется производительность различных приложений. Проведенные нами измерения относятся ко всем уровням и были выполнены в среде OpenSUSE 12.3 и Nvidia CUDA-5.5 на сервере с Intel Xeon E3-1240 и GPU K20c. Использовался также компилятор PGI Accelerator Fortran 13.10 с CUDA-расширениями.

На низшем уровне нами померена, в частности, пропускная способность обменов данными GPU-хост. Она составила около 6.4 Гбайт/с в каждом направлении. Наши более ранние данные для Nvidia C2050 с предыдущей архитектурой Fermi [1] дали около 6 Гбайт/с, что лишь немного ниже. В конечном счете это лимитируется пропускной способностью PCI-E (пиковая - 8 Гбайт/с), поэтому для приложений с большим обменом данными с хостом актуальна поддержка PCI-E v.3, имеющаяся в K40.

Для тестов среднего уровня отметим LU-декомпозицию на K40c, выполненную в пакете MAGMA [2]. На размерности матрицы 36K GPU обгоняет двухпроцессорный сервер с двумя 8-ядерными процессорами Xeon E5-2670/2.6 ГГц примерно в 4 раза, при размерности 2K быстродействие уже близко. На известной программе молекулярной динамики AMBER ускорение K20X относительно Intel Sandy Bridge около 8, на 80% быстрее Tesla M2090 [3].

Наши измерения для верхнего уровня были проведены для Nvidia Tesla C2050 и K20c по программе, реализующей квантовохимический метод PDM прямого построения матрицы плотности по фокиану в ортогональном базисе. Программа написана на Fortran-9X с CUDA-расширениями для создаваемого нами быстродействующего приближенного метода DFT (может использоваться и с полуэмпирическими методами типа AM/1) и для сверхбольших молекул обеспечивает линейное масштабирование времени расчета с размером системы при использовании технологии разреженных матриц с блочным портретом. Для молекулы полиглицина с базисом размерностью около 2.8 тысяч орбиталей ускорение на K20c по сравнению с C2050 составило около 2 раз (зависит от размера блока).

Время обменов данными с хостом почти не поменялось относительно C2050 и для блоков размерностью 200 составило 40% общего времени расчета. Это говорит о важности применения в GPU Kepler нового стандарта PCI-E v.3.

Работа поддержана РФФИ, проект 11-07-00470а.

### Литература

1. Кузьминский М. GPU для HPC — время пришло//Открытые системы. 2011.

№6. C.11-14

2. *Dongarra J., Dong T., Gates M., Haidar A., Tomov S., Yamazaki I.* MAGMA: a New Generation of Linear Algebra Libraries for GPU and Multicore Architectures [http://icl.utk.edu/projectsfiles/magma/pubs/MAGMA\\_1.4.pdf](http://icl.utk.edu/projectsfiles/magma/pubs/MAGMA_1.4.pdf)
3. NVIDIA Tesla® K20-K20X GPU Accelerators Benchmarks. Application Performance Technical Brief., Nvidia, Nov. 2012