

Программирование GPU на Фортран, Java, C#

Романенко А.А.

arom@ccfit.nsu.ru

Новосибирский государственный университет

Использование Фортран

- * Директивы в исходном коде
 - * OpenACC
- * Компиляторы, поддерживающие CUDA Fortran
 - * PGI accelerated fortran
 - * CRAY
- * Линковка модулей на CUDA C к программам на Фортран

OpenACC

- * Разработчики: **NVIDIA, PGI, CRAY, CAPS**
- * Директивы компилятору. Аналог OpenMP
- * Наличие API
- * Поддержка языков C/Fortran
- * Ресурсы:
 - * <http://www.openacc-standard.org>
 - * Руководство пользователя к PGI

OpenACC: модель выполнения

* Host

- * Выполняет большую часть кода;
- * Выделяет и освобождает память на GPU;
- * Управляет передачей данных и кода на GPU;
- * Управляет запуском ядер, параллельных циклов и синхронизацией;
- * Загружает результаты из памяти GPU;

* GPU

- * Выполняет ядра одно за другим
- * Генерирует синхронную/асинхронную передачу данных между хостом и GPU

OpenACC: директивы

* Fortran

```
!$acc directive [clause [, clause] ...]  
    structured block  
!$acc end directive
```

* C

```
#pragma acc directive [clause [, clause] ...]  
    structured block
```

* Компиляция программы

```
pgfortran -acc -Minfo=accel -ta=nvidia <file_name>  
pgcc -acc -Minfo=accel -ta=nvidia <file_name>
```

OpenACC: директивы

* Директивы

- * parallel
- * loop
- * kernels
- * и пр.

* Условия выполнения

- * if(condition)
- * async[(exp)]
- * num_gangs(exp)
- * num_workers(exp)
- * vector_length(exp)
- * reduction(operator:list)
- * и пр.

* Условия на данные

- * copy(list)
- * copyin(list)
- * copyout(list)
- * create(list)
- * present(list)
- * present_or_copy(list)
- * present_or_copyin(list)
- * present_or_copyout(list)
- * present_or_create(list)
- * deviceptr(list)
- * private(list)
- * firstprivate(list)

OpenACC. Ссылки

- * Стандарт

- * http://www.openacc.org/sites/default/files/OpenACC.1.0_0.pdf

- * Подсказки/рекомендации

- * <http://www.nvidia.com/docs/IO/117377/directives-tips-for-fortran.pdf>

- * <http://www.nvidia.com/docs/IO/117377/directives-tips-for-c.pdf>

- * Quick Reference card

- * http://www.openacc.org/sites/default/files/OpenACC_API_QuickRefGuide.pdf

CUDA Fortran

- * Отображение CUDA C на Фортран
- * Поддержка ряда операций
 - * средствами языка Фортран
 - * библиотечными функциями (Runtime API)
 - * **use cudafor**

CUDA Fortran: выделение памяти

- * Описание переменных

```
real, device, allocatable :: foo(:)
real, allocatable :: bar(:)
    attributes (device) :: bar
```

- * Выделение/освобождение памяти

```
allocate( foo(1:n), bar )
deallocate( foo )
err = cudaMalloc( bar, n )
err = cudaFree( bar )
```

CUDA Fortran: передача данных

```
real, device, allocatable :: da(:)
real, allocatable :: ha(:)
integer :: n
...
da(1:n) = ha(1:n)
...
err = cudaMemcpy(ha, da, n)
```

CUDA Fortran: запуск ядра

```
type(dim3) :: grid, block
...
grid = dim3(256, 1, 1)
block = dim3(512, 1, 1)
...
call kernel<<<grid, block>>>( параметры )
```

Запуск ядра асинхронный!

CUDA Fortran: ядро

```
attributes(global) subroutine cuj ( a, newa, n, m, w0, w1, w2, cc )
  real, value :: w0, w1, w2
  ...
  real, shared :: reduce(256)
  j = (blockidx%y-1)*blockdim%y + threadidx%y + 1
  i = (blockidx%x-1)*blockdim%x + threadidx%x + 1

  if( j < n .and. i < m )then
    newa(i,j) = w0 * a(i,j) + &
      w1 * (a(i-1,j) + a(i,j-1) + a(i+1,j) + a(i,j+1) ) + &
      w2 * (a(i-1,j-1) + a(i-1,j+1) + a(i+1,j-1) + a(i+1,j+1) )
    mychange = max( mychange, abs( newa(i,j) - a(i,j) ) )
  endif
  ir = (threadidx%y-1) * blockDim%x + threadidx%x
  reduce(ir) = mychange
  call syncthreads()
  ...
end subroutine
```

CUDA Fortran

- * **Компиляция**

- * `pgfortran -fast -O3 -Minfo=all
-Mcuda=cc13 -ta=nvidia:4.0 -o test test.f90`

- * **Ресурсы**

- * **PGI CUDA Fortran Compiler**
[<http://www.pgroup.com/resources/cudafortran.htm>]

CUDA C/Fortran vs OpenACC

- * **CUDA C/Fortran:**

- * + Высокая производительность при ручной настройке ядер;
- * + Инкрементальная переносимость на GPU;
- * - Только CUDA-платформы, несовместимость с другими;
- * - Необходима поддержка 2 наборов кода.

- * **OpenACC:**

- * + Возможна высокая производительность;
- * + Инкрементальная переносимость на GPU;
- * + Совместимость с другими не-CUDA-платформами;
- * + Необходима поддержка только 1 набора исходного кода;
- * - Сложный контроль над работой компилятора;
- * - Поддерживается только платными компиляторами.

Вызов CUDA ядра

- * Запуск CUDA-C-ядра из CUDA Fortran

```
interface
  attributes(global) subroutine saxpy(a,x,y,n) bind(c)
    real, device :: x(*), y(*)
    real, value  :: a
    integer, value :: n
  end subroutine
end interface
call saxpy<<<grid,block>>>(aa,xx,yy,nn)
```

- * Запуск CUDA-Fortran-ядра из CUDA C

```
extern __global__ void saxpy_( float a, float* x, float* y, int n );
saxpy_<<<grid,block>>>( a, x, y, n );
```

```
attributes(global) subroutine saxpy(a,x,y,n)
  real, value :: a
  real :: x(*), y(*)
  integer, value :: n
```

Вызов CUDA ядер из Фортран

CUDA C:

```
__global__ kernel (аргументы) {  
    ....  
}  
void some_function_(аргументы) {  
    ....  
    kernel <<<GS, BS>>> (аргументы);  
    ....  
}
```

Fortran:

```
....  
call some_function (аргументы);  
....
```


Программирование GPU на Java

- * AMD Aparapi

- * Только графические карты от AMD

- * Java код → OpenCL

- * <http://developer.amd.com/zones/java/aparapi/Pages/default.aspx>

- * java-gpu

- * Только графические карты с поддержкой CUDA

- * <http://code.google.com/p/java-gpu/>

Арағарі (пример)

```
final float inA[] = .... // get a float array of data from somewhere
final float inB[] = .... // get a float array of data from somewhere
                        // (inA.length==inB.length)
final float result = new float[inA.length];

for (int i=0; i<array.length; i++){
    result[i]=inA[i]+inB[i];
}
```

```
Kernel kernel = new Kernel(){
    @Override public void run(){
        int i= getGlobalId();
        result[i]=inA[i]+inB[i];
    }
};
Range range = Range.create(result.length);
kernel.execute(range);
```

java-gpu (пример)

```
@Parallel(loops = {"y", "x"})
public void compute() {
    for(int y = 0; y < height; y++) {
        for(int x = 0; x < width; x++) {
            float Zr = 0.0f;
            float Zi = 0.0f;
            ...
            data[y][x] = (short)((i * 255) / iterations);
        }
    }
}
```

```
java -jar Parallel.jar samples.Mandelbrot
```

```
java -cp ../Parallel.jar samples.Mandelbrot 2000 2000 out.png
```

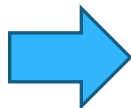
Программирование GPU на Java

- * [jcuda.org](http://www.jcuda.org/)
 - * Обертки для CUDA runtime API и driver API
 - * Ядра на CUDA C
 - * <http://www.jcuda.org/>
- * [jocl.org](http://www.jocl.org/)
 - * Обертка для OpenCL API
 - * <http://www.jocl.org/>
- * Lightweight Java Game Library (LWJGL)
 - * Обертка для OpenCL API
 - * <http://www.lwjgl.org>
- * JavaCL
 - * Обертка для OpenCL API
 - * Поддержка видеокарт от AMD и NVida
 - * <http://code.google.com/p/javacl/>
- * пр.

```
import java.lang.*;
public class Hello {
    static { System.loadLibrary("Hello"); }
    public static native String getMessage();
    public static void main( String[] args ) {
        System.out.println( getMessage() );
        System.exit(0);
    }
}
```



```
javac Hello.java
javah Hello
```



```
//Hello.h
#include <jni.h>
#ifndef _Included_Hello
#define _Included_Hello
#ifdef __cplusplus
extern "C" {
#endif
JNIEXPORT jstring JNICALL
    Java_Hello_getMessage (JNIEnv *, jclass);
#ifdef __cplusplus
}
#endif
#endif
```

<http://xyplot.com/jni.simple.htm>

Программирование GPU на C#

- * CUDA.NET
 - * Обертка к CUDA runtime API
 - * Страница проекта удалена
- * GPU.NET
 - * Разметка кода на C#
 - * Только устройства с поддержкой CUDA
 - * Пробная лицензия на 30 дней
 - * <http://www.tidepowerd.com/gpu-net>
- * Accelerator
 - * Исследовательский проект от Microsoft
 - * Последнее обновление в 2010 году
 - * <http://research.microsoft.com/en-us/projects/Accelerator/>
- * Cloo
 - * Обертка к OpenCL
 - * <http://cloo.sourceforge.net/>

DLL:

```
#include <thrust/device_vector.h>
#include <thrust/sort.h>
#include <thrust/copy.h>
#include <thrust/detail/type_traits.h>

extern "C" __declspec(dllexport) void __cdecl GPUSort(int*, unsigned int);
extern void GPUSort(int* data, unsigned int numElements) {
    thrust::device_vector<int> d_data(data, data + numElements);
    thrust::stable_sort(d_data.begin(), d_data.end());
    thrust::copy(d_data.begin(), d_data.end(), data);
}
```

C#:

```
[DllImport("GPUSort.dll", CallingConvention = CallingConvention.Cdecl)]
public static extern void GPUSort(
    [MarshalAsAttribute(UnmanagedType.LPArray,
        ArraySubType = UnmanagedType.I4)]
    int[] data, uint numElements);
```