# Rock' em Graphic Cards

Agnes Meyder

27.12.2013, 16:00

# Layout

# Layout

# Big Data

Buzz word for: "Pattern Recognition in large Datasets"

NASA, ESA, and M. Livio and the Hubble 20th Anniversary Team (STScI) • Hubble Space Telescope WFC3/UVIS • STScI-PRC10-13e
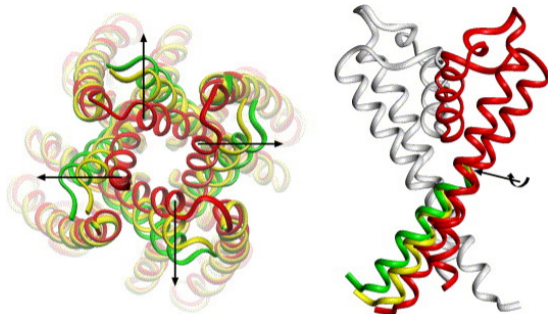
Figure: hERG - Possible Opening Movement[1]

# Sponsors for my Summer School

# Layout

# Parallelism

**data-based**                **task-based**

explicit

CUDA    OpenCL

C++AMP                OpenMPI

OpenMP    OpenACC

abstract

# Parallelism

| data-based | task-based |
| --- | --- |

# Types of Parallelism

**Data-Parallelism**

- Fry omelets for 8.000 people.

**Task-Parallelism**

- Cook a 5-course menu for one person.

Cook a menu for 8.000 people with an omelet as part of the dessert.

# PP - Possible Problems

- ▶ Kitchen too small
- ▶ Only apprentices available
- ▶ Not enough frying pans
- ▶ Delivery paths too small (only one can access the fridge)
- ▶ Only one can write a new recipe into the book
- ▶ Transport 30 eggs in one go
- ▶ Serve the courses in the correct order

# PP - Possible Problems

| | | |
|---|---|---|
| Kitchen too small | $\Rightarrow$ | Global capacity limitation |
| Only apprentices available | $\Rightarrow$ | Processor complexity limited |
| Not enough frying pans | $\Rightarrow$ | Concurrent task number limited |
| Delivery paths too small (only one can access the fridge) | $\Rightarrow$ | Band width limitation |
| Only one can write a new recipe into the book | $\Rightarrow$ | Read-write access limitation |
| Transport 30 eggs in one go | $\Rightarrow$ | coalescing memory access |
| Serve the courses in the correct order | $\Rightarrow$ | Synchronization issues |

# Layout

# Parallelism - Old Standards

|  | **data-based** | **task-based** |
|---|---|---|

explicit

OpenMPI

OpenMP

implicit

# OpenMPI - Installation Sprint

OpenMPI is one implementation of the standard MPI.

- ▶ Install openmpi and its devel packages
- ▶ Switch linker and compiler to mpicc
- ▶ Compile: `mpicc -I/usr/lib64/mpi/gcc/openmpi/include ...`
- ▶ Link: `mpicc -o "openmpi" ./src/openmpi.o -lmpi`
- ▶ Run code with: `mpirun -np 1 -hostfile hostfile ./openmpi`

**hostfile:** a file with the names of usable hosts. "localhost" is fitting for testing purposes.
**Dcoumentation**: `http://www.open-mpi.org/`

# OpenMPI - "Hello World"[2]

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <mpi.h>

5  int main(int argc, char *argv[]) {
6    int numprocs, rank, namelen;
7    char processor_name[MPI_MAX_PROCESSOR_NAME];

9    MPI_Init(&argc, &argv);
10   MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
11   MPI_Comm_rank(MPI_COMM_WORLD, &rank);
12   MPI_Get_processor_name(processor_name, &namelen);

14   printf("Process %d on %s out of %d\n",
15   rank, processor_name, numprocs);

17   //cook one course per task ?

19   MPI_Finalize();
20 }
```
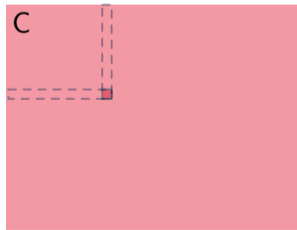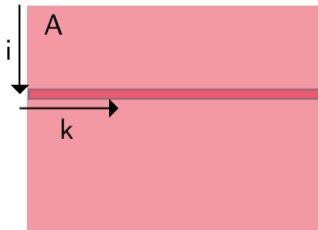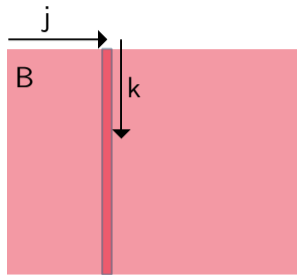
---

[2] http://www.linux-mag.com/id/5759/

# Typical (Parallel) Task: Matrix-Multiplication

# Matrix-Multiplication[3]

$$C_{ij} = \sum_{k=1}^{m} A_{ik} \cdot B_{kj}$$

# Matrix-Multiplication

```
1    const int NUM_ROWS_A  = 200;
2    const int NUM_COLS_A  = 300;
3    const int NUM_ROWS_B  = NUM_COLS_A;
4    const int NUM_COLS_B  = 400;

6    int *A = fillMatrix(NUM_ROWS_A, NUM_COLS_A);
7    int *B = fillMatrix(NUM_ROWS_B, NUM_COLS_B);
8    int *C =  (int*)malloc(sizeof(int)*NUM_ROWS_A  * NUM_COLS_B);

10   for (unsigned int i=0; i < NUM_ROWS_A; i++)
11   {
12     for(unsigned int  j=0; j != NUM_COLS_B; ++j)
13     {
14       C[i*NUM_COLS_B +j] = 0;
15       for(unsigned int  k=0; k != NUM_COLS_A; ++k)
16       {
17         C[i*NUM_COLS_B +j] += A[i*NUM_COLS_A +k ]
18           * B[k*NUM_COLS_B +j ];
19       }
20     }
21   }
22   // Do what you need with C ...
23   free(A); free(B); free(C);
```

# OpenMP - "Hello World" [4]

```c
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <omp.h>

5  int main(void) {
6    //everything in the bracket is parallel!
7    #pragma omp parallel
8    {
9    printf("Hello World !! This program has %d threads.. "
10     "This line is printed by Thread: %d \n",
11     omp_get_num_threads(),
12     omp_get_thread_num());
13   }
14   return EXIT_SUCCESS;
15 }
```

---

[4]http://www.linux-mag.com/id/5759/

# OpenMP - Matrix-Multiplication I

```c
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <omp.h>

5  int * fillMatrix(const int num_rows, const int num_cols)
6  {
7    int i;
8    int *pointer = (int*)malloc(sizeof(int)*num_cols * num_rows);
9    for(i = 0; i < num_cols * num_rows; i++) {
10     pointer[i] = i;
11   }
12   return pointer;
13 }

15 int main(void) {
16   const int NUM_ROWS_A = 2;
17   const int NUM_COLS_A = 3;
18   const int NUM_ROWS_B = NUM_COLS_A ;
19   const int NUM_COLS_B = 4;

21   int *A = fillMatrix(NUM_ROWS_A, NUM_COLS_A);
22   int *B = fillMatrix(NUM_ROWS_B, NUM_COLS_B);
23   int *C =  (int*)malloc(sizeof(int)*NUM_ROWS_A  * NUM_COLS_B);
```

# OpenMP - Matrix-Multiplication II

```c
1    //everything in the bracket is parallel!
2    #pragma omp parallel
3    {
4      printf("Hello World! This program has %d threads.. "
5          "This line is printed by Thread: %d \n",
6          omp_get_num_threads(),
7          omp_get_thread_num());
8      for(unsigned int i = 0; i != NUM_ROWS_A; ++i)
9      {
10       for(unsigned int j=0; j != NUM_COLS_B; ++j)
11       {
12         C[i*NUM_COLS_B +j] = 0;
13         for(unsigned int k=0; k != NUM_COLS_A; ++k)
14         {
15           C[i*NUM_COLS_B +j] += A[i*NUM_COLS_A +k ]
16             * B[k*NUM_COLS_A +j ];
17           printf( "C[%d, %d] is calculated by Thread: %d \n",
18               i, j, omp_get_thread_num());
19         }
20       }
21     }
22   } // end pragma
```

# OpenMP - Matrix-Multiplication III

```c
1     for(unsigned int i = 0; i != NUM_ROWS_A; ++i)
2     {
3       for(unsigned int j=0; j != NUM_COLS_B; ++j)
4       {
5         printf( "Row: %d Column: %d Value: %d\n" , i,  j,
6           C[i*NUM_COLS_B +j]);
7       }
8     }

10    free(A);
11    free(B);
12    free(C);
13    return EXIT_SUCCESS;
14  }
```

# OpenMP MM- Output

```
1    Hello World! This program has 4 threads..
2       This line is printed by Thread: 2
3  C [0, 0] is calculated by Thread: 2
4  ....
5  Hello World! This program has 4 threads..
6       This line is printed by Thread: 3
7  C [0, 0] is calculated by Thread: 3
8  ...
9  Hello World! This program has 4 threads..
10      This line is printed by Thread: 0
11 Hello World! This program has 4 threads..
12      This line is printed by Thread: 1
13 C [0, 0] is calculated by Thread: 1
14 ...
```

# OpenMP - Matrix-Multiplication IV

```
1    int i, j, k, chunk;
2  ...
3    const int CHUNKSIZE = 4;
4  ...
5    chunk = CHUNKSIZE;

7    #pragma omp parallel shared(A,B,C,chunk, i) private(j, k)
8    {
9      printf("Numer of launched threads: %d\n",
10       omp_get_num_threads());
11     #pragma omp for schedule(dynamic,chunk) nowait
12     for (i=0; i < NUM_ROWS_A; i++)
13     {
14       for(j=0; j != NUM_COLS_B; ++j)
15       {
16         C[i*NUM_COLS_B +j] = 0;
17         for(k=0; k != NUM_COLS_A; ++k)
18         {
19           C[i*NUM_COLS_B +j] += A[i*NUM_COLS_A +k ]
20             * B[k*NUM_COLS_B +j ];
21         }
22       }
23     } /* end of parallel for loop */
24   } /* end of parallel section */
```

# OpenMP - Installation Sprint

- Install libgomp
- Compile: `gcc  -fopenmp ...`
- Link: `gcc -o "openmp"  ./src/openmp.o   -lgomp`

**Documentation:** `http://www.openmp.org/mp-documents/OpenMP4.0.0.pdf`

# Layout
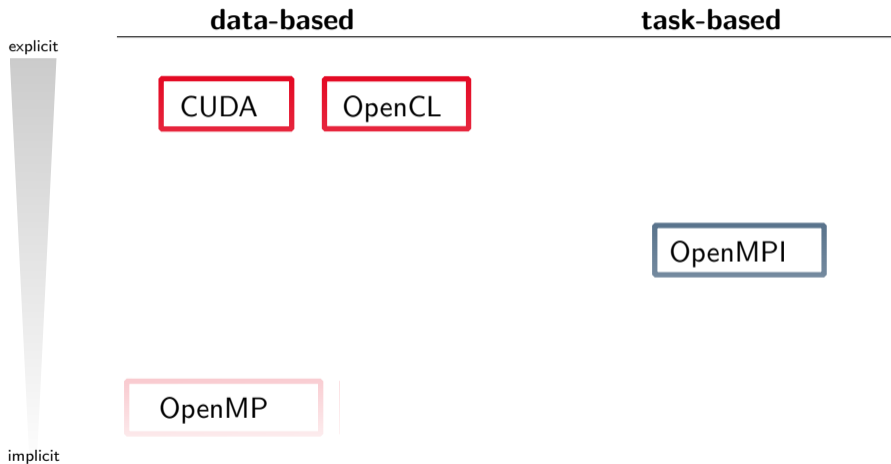
# Parallelism

# Accelerator Terminology

Some words:

- Host = master thread, often CPU
- Device = helper threads, often on accelerator cards
- Kernel = functions, running on the device

Some words:

- In order to run OpenCL- or CUDA-Code, please install the proprietary drivers for your graphic card!

# CUDA- Matrix Multiplication I

```c
1  #include <stdio.h>
2  #include <stdlib.h>

4  __global__ void matrixMultiply(int *A, int *B, int *C,
5      int NUM_ROWS_A, int NUM_COLS_A,
6      int NUM_ROWS_B, int NUM_COLS_B,
7      int NUM_ROWS_C, int NUM_COLS_C) {
8    int row = blockIdx.x * blockDim.x + threadIdx.y;
9    int col = blockIdx.x * blockDim.x + threadIdx.x;

11   if((row < NUM_ROWS_C) && (col < NUM_COLS_C))
12   {
13     int sum = 0;
14     for (int k = 0; k < NUM_COLS_A; ++k)
15     {
16       sum += A[row*NUM_COLS_A + k] * B[k*NUM_COLS_B + col];
17     }
18     C[row*NUM_COLS_C+col] = sum;
19   }

21 }

23 int main(void) {
```

# CUDA - Matrix-Multiplication II

```
1    int *A = fillMatrix(NUM_ROWS_A, NUM_COLS_A);
2    int *B = fillMatrix(NUM_ROWS_B, NUM_COLS_B);
3    int *C = (int*)malloc(sizeof(int)*NUM_ROWS_C * NUM_COLS_C);
4    for(int i = 0; i != NUM_ROWS_C * NUM_COLS_C; ++i)
5    {
6      C[i] = -1;
7    }

9    cudaMalloc((void**) &dA, NUM_ROWS_A*NUM_COLS_A*sizeof(int));
10   cudaMalloc((void**) &dB, NUM_ROWS_B*NUM_COLS_B*sizeof(int));
11   cudaMalloc((void**) &dC, NUM_ROWS_C*NUM_COLS_C*sizeof(int));

13   cudaMemcpy(dA, A, NUM_ROWS_A*NUM_COLS_A*sizeof(int), cudaMemcpyHostToDevice);
14   cudaMemcpy(dB, B, NUM_ROWS_B*NUM_COLS_B*sizeof(int), cudaMemcpyHostToDevice);

16   // Initialize the grid and block dimensions here
17   dim3 dimGrid(1, 1,1);
18   dim3 dimBlock(NUM_COLS_C, NUM_ROWS_C, 1);

20   // Launch the GPU Kernel here
21   matrixMultiply<<<dimGrid, dimBlock>>>(dA, dB, dC, NUM_ROWS_A, NUM_COLS_A,
22       NUM_ROWS_B, NUM_COLS_B, NUM_ROWS_C, NUM_COLS_C);
```

# CUDA - Matrix-Multiplication III

```
1    cudaThreadSynchronize();

3    // Copy the GPU memory back to the CPU here
4    cudaMemcpy(C, dC, NUM_ROWS_C*NUM_COLS_C*sizeof(int), cudaMemcpyDeviceToHost);

6    cudaFree(dC);
7    cudaFree(dB);
8    cudaFree(dA);

10   // Do what you need with C ...

12   free(A);
13   free(B);
14   free(C);
15   return EXIT_SUCCESS;
16   }
```

# CUDA- Matrix Multiplication I

```
1  __global__ void matrixMultiply(int *A, int *B, int *C, int NUM_ROWS_A, int NUM_CO
2      int NUM_ROWS_B, int NUM_COLS_B, int NUM_ROWS_C, int NUM_COLS_C) {
3    int row = blockIdx.x * blockDim.x + threadIdx.y;
4    int col = blockIdx.x * blockDim.x + threadIdx.x;

6    if((row < NUM_ROWS_C) && (col < NUM_COLS_C))
7    {
8      int sum = 0;
9      for (int k = 0; k < NUM_COLS_A; ++k)
10     {
11       sum += A[row*NUM_COLS_A + k] * B[k*NUM_COLS_B + col];
12     }
13     C[row*NUM_COLS_C+col] = sum;
14   }
15 }

17 int main(void){

19   dim3 dimGrid(1, 1,1);
20   dim3 dimBlock(NUM_COLS_C, NUM_ROWS_C, 1);

22   matrixMultiply<<<dimGrid, dimBlock>>>(dA, dB, dC, NUM_ROWS_A, NUM_COLS_A,
23       NUM_ROWS_B, NUM_COLS_B, NUM_ROWS_C, NUM_COLS_C);
24   cudaThreadSynchronize();
```

# CUDA - Installation Sprint

- Install CUDA SDK and your os CUDA package (check compatibility of gcc!)
- Compile: `nvcc ... --compiler-bindir /path/to/special/gcc`
- Link: `gcc -L/usr/local/cuda/lib64 ... -lcudart?`

**Documentation:** `http://docs.nvidia.com/cuda/`

# OpenCL - Vector Addition I[5]

```
1   #include <stdio.h>
2   #include <stdlib.h>
3   #ifdef __APPLE__
4   #include <OpenCL/opencl.h>
5   #else
6   #include <CL/cl.h>
7   #endif

9   #define MAX_SOURCE_SIZE (0x100000)

11  int main(void) {
12    // Create the two input vectors
13    int i;
14    const int LIST_SIZE = 1024;
15    int *A = (int*)malloc(sizeof(int)*LIST_SIZE);
16    int *B = (int*)malloc(sizeof(int)*LIST_SIZE);
17    for(i = 0; i < LIST_SIZE; i++) {
18      A[i] = i;
19      B[i] = LIST_SIZE - i;
20    }
```

[5]http://www.thebigblob.com/getting-started-with-opencl-and-gpu-computing/

# OpenCL - Vector Addition II

```
1    // Load the kernel source code into the array source_str
2    FILE *fp;
3    char *source_str;
4    size_t source_size;
5    fp = fopen("src/vector_add_kernel.cl", "r");
6    if (!fp) {
7      fprintf(stderr, "Failed to load kernel.\n");
8      exit(1);
9    }
10   source_str = (char*)malloc(MAX_SOURCE_SIZE);
11   source_size = fread( source_str, 1, MAX_SOURCE_SIZE, fp);
12   fclose( fp );

14   // Get platform and device information
15   cl_platform_id platform_id = NULL;
16   cl_device_id device_id = NULL;
17   cl_uint ret_num_devices;
18   cl_uint ret_num_platforms;
19   cl_int ret = clGetPlatformIDs(1, &platform_id, &ret_num_platforms);
20   ret = clGetDeviceIDs( platform_id, CL_DEVICE_TYPE_GPU, 1,
21       &device_id, &ret_num_devices);
```

# OpenCL - Vector Addition III

```
1    // Create an OpenCL context
2    cl_context context = clCreateContext( NULL, 1, &device_id,
3      NULL, NULL, &ret);
4    // Create a command queue
5    cl_command_queue command_queue = clCreateCommandQueue(
6      context, device_id, 0, &ret);
7    // Create memory buffers on the device for each vector
8    cl_mem a_mem_obj = clCreateBuffer(context, CL_MEM_READ_ONLY,
9      LIST_SIZE * sizeof(int), NULL, &ret);
10   cl_mem b_mem_obj = clCreateBuffer(context, CL_MEM_READ_ONLY,
11     LIST_SIZE * sizeof(int), NULL, &ret);
12   cl_mem c_mem_obj = clCreateBuffer(context, CL_MEM_WRITE_ONLY,
13     LIST_SIZE * sizeof(int), NULL, &ret);
14   // Copy the lists A and B to their respective memory buffers
15   ret = clEnqueueWriteBuffer(command_queue, a_mem_obj,
16     CL_TRUE, 0, LIST_SIZE * sizeof(int), A, 0, NULL, NULL);
17   ret = clEnqueueWriteBuffer(command_queue, b_mem_obj,
18     CL_TRUE, 0, LIST_SIZE * sizeof(int), B, 0, NULL, NULL);
19   // Create a program from the kernel source
20   cl_program program = clCreateProgramWithSource(context, 1,
21     (const char **)&source_str, (const size_t *)&source_size,
22     &ret);
```

# OpenCL - Vector Addition IV

```
1    // Build the program
2    ret = clBuildProgram(program, 1, &device_id, NULL,
3      NULL, NULL);

5    // Create the OpenCL kernel
6    cl_kernel kernel = clCreateKernel(program,
7      "vector_add", &ret);

9    // Set the arguments of the kernel
10   ret = clSetKernelArg(kernel, 0, sizeof(cl_mem),
11     (void *)&a_mem_obj);
12   ret = clSetKernelArg(kernel, 1, sizeof(cl_mem),
13     (void *)&b_mem_obj);
14   ret = clSetKernelArg(kernel, 2, sizeof(cl_mem),
15     (void *)&c_mem_obj);

17   // Execute the OpenCL kernel on the list
18   size_t global_item_size = LIST_SIZE; // Process the entire lists
19   size_t local_item_size = 64; // Process in groups of 64
20   ret = clEnqueueNDRangeKernel(command_queue, kernel, 1,
21     NULL, &global_item_size, &local_item_size, 0, NULL, NULL);
```

# OpenCL - Vector Addition V[6]

```
1      // Read the memory buffer C on the device to the local variable C
2      int *C = (int*)malloc(sizeof(int)*LIST_SIZE);
3      ret = clEnqueueReadBuffer(command_queue, c_mem_obj,
4        CL_TRUE, 0, LIST_SIZE * sizeof(int), C, 0, NULL, NULL);
5      // Display the result to the screen
6      for(i = 0; i < LIST_SIZE; i++)
7        printf("%d + %d = %d\n", A[i], B[i], C[i]);
8      // Clean up
9      ret = clFlush(command_queue);
10     ret = clFinish(command_queue);
11     ret = clReleaseKernel(kernel);
12     ret = clReleaseProgram(program);
13     ret = clReleaseMemObject(a_mem_obj);
14     ret = clReleaseMemObject(b_mem_obj);
15     ret = clReleaseMemObject(c_mem_obj);
16     ret = clReleaseCommandQueue(command_queue);
17     ret = clReleaseContext(context);
18     free(A);
19     free(B);
20     free(C);
21     return 0;
22   }
```

---

# OpenCL - Vector Addition Summary

```
1    // Get platform and device information
2    cl_platform_id platform_id = NULL;
3    cl_device_id device_id = NULL;
4    cl_uint ret_num_devices;
5    cl_uint ret_num_platforms;
6    cl_int ret = clGetPlatformIDs(1, &platform_id, &ret_num_platforms);
7    ret = clGetDeviceIDs( platform_id, CL_DEVICE_TYPE_GPU, 1,
8        &device_id, &ret_num_devices);
9  ...
10   // Execute the OpenCL kernel on the list
11   size_t global_item_size = LIST_SIZE; // Process the entire lists
12   size_t local_item_size = 64; // Process in groups of 64
13   ret = clEnqueueNDRangeKernel(command_queue, kernel, 1,
14     NULL, &global_item_size, &local_item_size, 0, NULL, NULL);
15 ...
16    // Read the memory buffer C on the device to the local variable C
17   int *C = (int*)malloc(sizeof(int)*LIST_SIZE);
18   ret = clEnqueueReadBuffer(command_queue, c_mem_obj,
19     CL_TRUE, 0, LIST_SIZE * sizeof(int), C, 0, NULL, NULL);
```

# OpenCL - Installation Sprint

- ▶ Install Khronos OpenCL-Headers and if necessary SDK of AMD/NVIDIA
- ▶ Compile: `gcc ...`
- ▶ Link: `gcc ... -lOpenCL`

**Documentation:**
http://www.khronos.org/opencl/
http://www.khronos.org/files/opencl-1-2-quick-reference-card.pdf

# General Code Structure

- ▶ Define accelerator device
- ▶ Initialize: memory on host and device
- ▶ Specify amount and dimensions of necessary kernel launches
- ▶ Launch Kernel
- ▶ Cleanup: Store result, free data

# Parallelism

|  | **data-based** | **task-based** |
| --- | --- | --- |

explicit

CUDA    OpenCL

OpenMPI

OpenMP    OpenACC

implicit

# OpenMP - Matrix-Multiplication IV

```
1    #pragma omp parallel shared(A,B,C,chunk, i) private(j, k)
2    {
3      printf("Numer of launched threads: %d\n",
4        omp_get_num_threads());

6      #pragma omp for schedule(dynamic,chunk) nowait
7      for (i=0; i < NUM_ROWS_A; i++)
8      {

10       for(j=0; j != NUM_COLS_B; ++j)
11       {
12         C[i*NUM_COLS_B +j] = 0;
13         for(k=0; k != NUM_COLS_A; ++k)
14         {
15           C[i*NUM_COLS_B +j] += A[i*NUM_COLS_A +k ]
16             * B[k*NUM_COLS_B +j ];
17         }
18       }
19     } /* end of parallel for loop */
20    }  /* end of parallel section */
```

# OpenACC - Matrix-Multiplication

**OpenMP**:

```
1    #pragma omp parallel shared(A,B,C,chunk, i) private(j, k)
2    {
3      #pragma omp for schedule(dynamic,chunk) nowait
4      for (i=0; i < NUM_ROWS_A; i++)
5      { ...
```

**OpenACC**:

```
1    #pragma acc parallel
2    {
3      #pragma acc loop
4      for (i=0; i < NUM_ROWS_A; i++)
5      {

7        for(j=0; j != NUM_COLS_B; ++j)
8        {
9          C[i*NUM_COLS_B +j] = 0;
10         for(k=0; k != NUM_COLS_A; ++k) { ...
```

# OpenACC - Installation Sprint?

Already available in cray, cups, pgi compilers.

Available in **gcc** in 2015?!

# Tiled Matrix Multiplication I



$$C_{ij} = \sum_{k=1}^{m} A_{ik} \cdot B_{kj}$$

$$C_{ij} = \sum_{k=1}^{m} A_{ik} \cdot B_{kj}$$

# Tiled Matrix Multiplication - CUDA Warp Divergence



$$C_{ij} = \sum_{k=1}^{m} A_{ik} \cdot B_{kj}$$

# NVIDIA[7] - Simplified Architecture Scheme
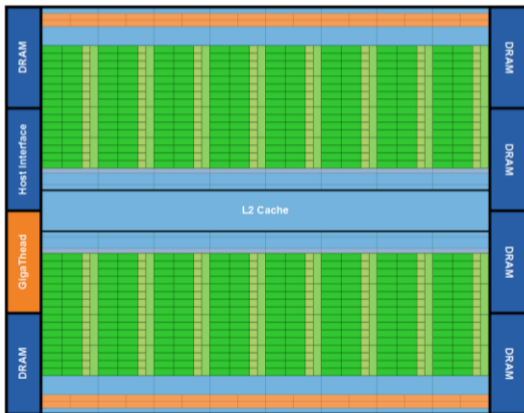
# Thread vs. Block vs. Grid Dimensions in CUDA

| CUDA | OpenCL | OpenACC |
|------|--------|---------|
| thread | work-item | vector |
| | | (worker) |
| block | work-group | worker / gang |
| grid | | (gang) |

# Hardware

- Graphic Cards (AMD, NVIDIA)
- Xeon Phi (Intel)
- parallela (`http://parallela.org`
- FPGAs (field programmable gate arrays)

Fermi's 16 SM are positioned around a common L2 cache. Each SM is a vertical rectangular strip that contain an orange portion (scheduler and dispatch), a green portion (execution units), and light blue portions (register file and L1 cache).

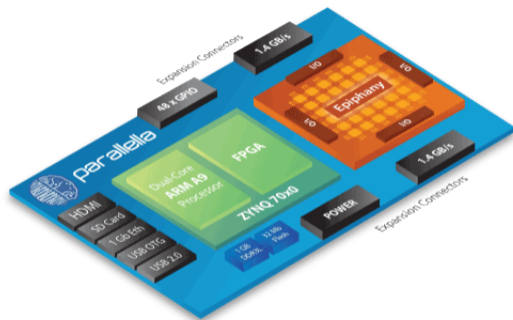[8]http://www.nvidia.com/content/PDF/fermi_white_papers/NVIDIA_Fermi_Compute_Architecture_Whitepaper.pdf

# Intel Xeon Phi[9] - 2000€

[9]http://software.intel.com/en-us/articles/
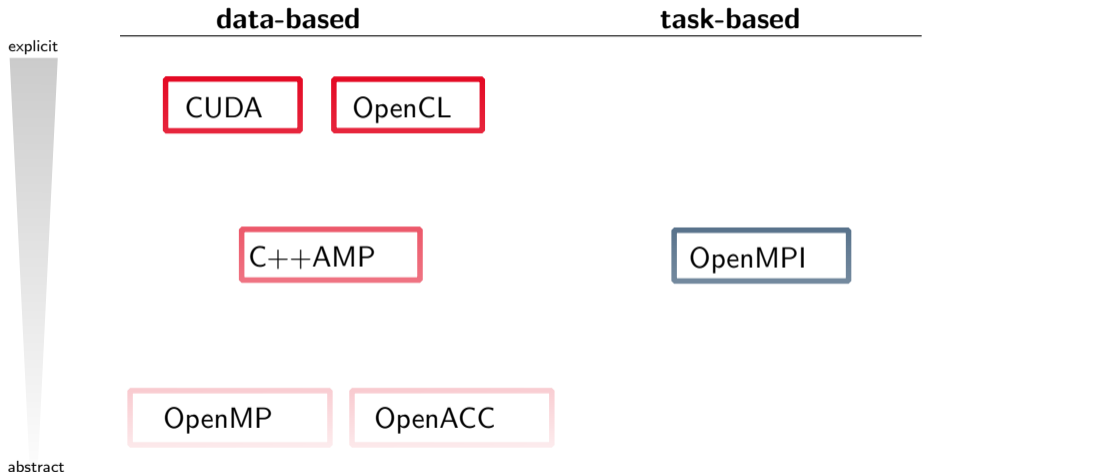intel-xeon-phi-coprocessor-codename-knights-corner

# Parallela[10] - $99

# Parallelism

# Layout

# I did not mention:

- POSIXThreads
- BOOST Threads
- Cilk (in C++: Cilk Plus)
- Intel Thread Building Blocks
- task-scheduling programs such as StarPU, oomps

… and languages such as Python with rudimentary multiprocessing abilities.

# Parallelism