

# An Introduction to OpenCL for Scientific Computing

Florian Rudolf<sup>1</sup>, Karl Rupp<sup>1,2</sup>, Josef Weinbub<sup>1</sup>

<http://karlrupp.net/>

<sup>1</sup>Institute for Microelectronics

<sup>2</sup>Institute for Analysis and Scientific Computing  
Technische Universität Wien, Austria



Workshop Programming of Heterogeneous Systems in Physics

July 14th, 2014

## Positions

PhD student at TU Wien (2009-2011)

Postdoc at Argonne Natl. Lab. (09/2012-09/2013)

Postdoc at TU Wien (09/2013-current)

## Research Interests

Semiconductor device simulation

Numerical solution of PDEs

Parallel computing

## Software Development

PETSc

ViennaCL

ViennaSHE

...



## Debunking the 100X GPU vs. CPU Myth: An Evaluation of Throughput Computing on CPU and GPU

Victor W Lee<sup>†</sup>, Changkyu Kim<sup>†</sup>, Jatin Chhugani<sup>†</sup>, Michael Deisher<sup>†</sup>,  
Daehyun Kim<sup>†</sup>, Anthony D. Nguyen<sup>†</sup>, Nadathur Satish<sup>†</sup>, Mikhail Smelyanskiy<sup>†</sup>,  
Srinivas Chennupati<sup>\*</sup>, Per Hammarlund<sup>\*</sup>, Ronak Singhal<sup>\*</sup> and Pradeep Dubey<sup>†</sup>

victor.w.lee@intel.com

<sup>†</sup>Throughput Computing Lab,  
Intel Corporation

<sup>\*</sup>Intel Architecture Group,  
Intel Corporation

### ABSTRACT

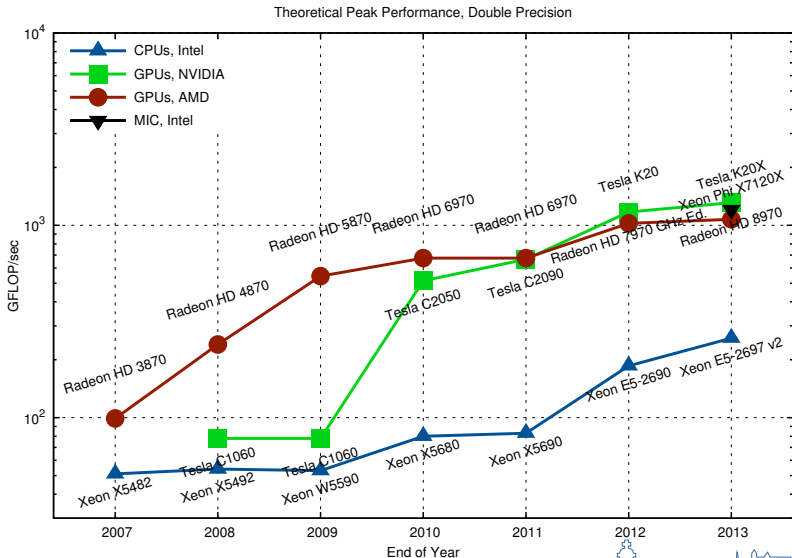
Recent advances in computing have led to an explosion in the amount of data being generated. Processing the ever-growing data in a timely manner has made throughput computing an important aspect for emerging applications. Our analysis of a set of important throughput computing kernels shows that there is an ample amount of parallelism in these kernels which makes them suitable for today's multi-core CPUs and GPUs. In the past few years there have been many studies claiming GPUs deliver substantial speedups (between 10X and 1000X) over multi-core CPUs on these kernels. To understand where such large performance difference comes from, we perform a rigorous performance analysis and find that after ap-

### 1. INTRODUCTION

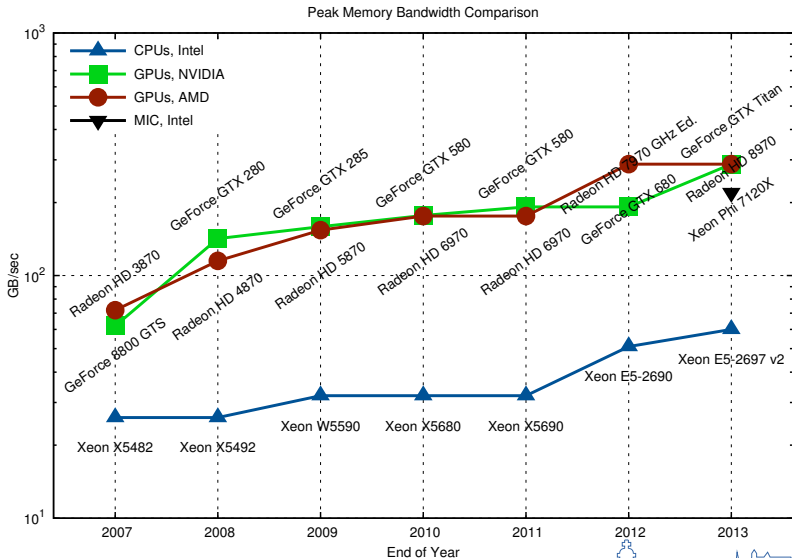
The past decade has seen a huge increase in digital content as more documents are being created in digital form than ever before. Moreover, the web has become the medium of choice for storing and delivering information such as stock market data, personal records, and news. Soon, the amount of digital data will exceed exabytes ( $10^{18}$ ) [31]. The massive amount of data makes storing, cataloging, processing, and retrieving information challenging. A new class of applications has emerged across different domains such as database, games, video, and finance that can process this huge amount of data to distill and deliver appropriate content to users. A distinguishing feature of these applications is that they



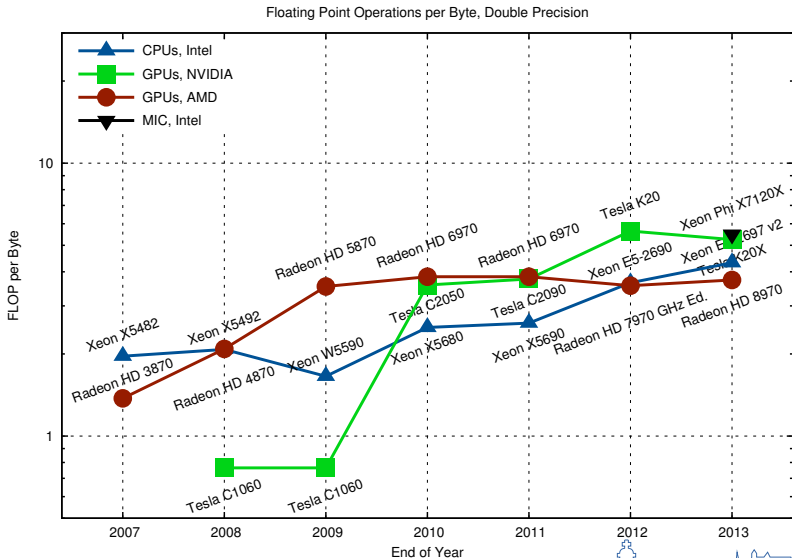
# 100x Speed-Up?



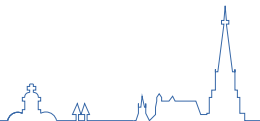
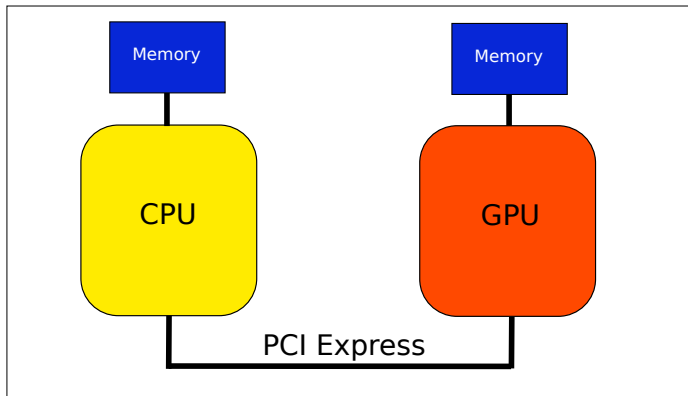
# 100x Speed-Up?



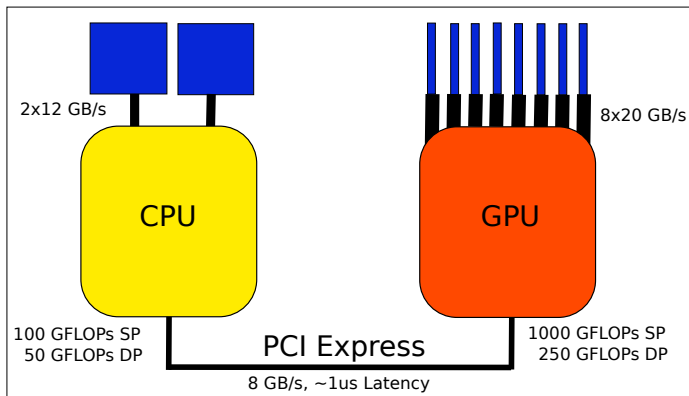
# 100x Speed-Up?



## Computing Architecture Schematic



## Computing Architecture Schematic



Good for large FLOP-intensive tasks, high memory bandwidth

PCI-Express can be a bottleneck

» 10-fold speedups (usually) not backed by hardware



# Table of Contents

100x Speedup!?

## About OpenCL

OpenCL Execution Model

OpenCL Kernel Language

Comparison with CUDA and OpenACC

Portable Performance

Summary



# History of OpenCL

## Prior to 2008

OpenCL developed by Apple Inc.

## 2008

OpenCL working group formed at Khronos Group  
OpenCL specification 1.0 released

## 2010

OpenCL 1.1 (multi-device, subbuffer manipulation)

## 2011

OpenCL 1.2 (device partitioning)

## 2013

OpenCL 2.0 (shared virtual memory, SPIR, etc.)



**OpenCL**



## Advantages

Not restricted to a single vendor: Intel, NVIDIA, AMD, ...

Just a shared C-library

Does not rely on compiler magic



## Advantages

Not restricted to a single vendor: Intel, NVIDIA, AMD, ...

Just a shared C-library

Does not rely on compiler magic

## Disadvantages

Not restricted to a single vendor

Boilerplate code required

Portable performance?



# Table of Contents

100x Speedup!?

About OpenCL

**OpenCL Execution Model**

OpenCL Kernel Language

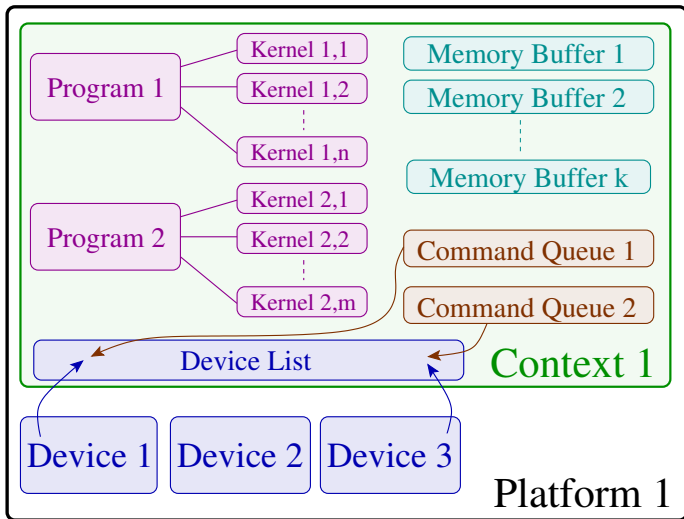
Comparison with CUDA and OpenACC

Portable Performance

Summary



# OpenCL Platform Model



<https://github.com/karlsruhp/phsp2014>



```
// Setup context and queue
context = clCreateContextFromType(NULL, CL_DEVICE_TYPE_GPU, NULL, NULL, NULL);
queue = clCreateCommandQueue(context, NULL, 0, NULL);

// Create memory buffers
memobjs[0] = clCreateBuffer(context, CL_MEM_READ_WRITE, sizeof(float)*2*num_entries
, NULL, NULL);
memobjs[1] = clCreateBuffer(context, CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
    sizeof(float)*2*num_entries, srcA, NULL);

// Create OpenCL program and kernels
program = clCreateProgramWithSource(context, 1, &kernel_src, NULL, NULL);
clBuildProgram(program, 0, NULL, NULL, NULL, NULL);

kernel = clCreateKernel(program, "my_kernel", NULL);
clSetKernelArg(kernel, 0, sizeof(cl_mem), (void *)&memobjs[0]);
clSetKernelArg(kernel, 1, sizeof(cl_mem), (void *)&memobjs[1]);
clSetKernelArg(kernel, 2, sizeof(float)*(local_work_size[0]+1)*16, NULL);

// Run an OpenCL kernel:
global_work_size[0] = 128;
    local_work_size[0] = 64;
clEnqueueNDRangeKernel(queue, kernel, 1, NULL, global_work_size, local_work_size,
    0, NULL, NULL);
```

## Issues

“Where is the error?”

Manage OpenCL handles





# Table of Contents

100x Speedup!?

About OpenCL

OpenCL Execution Model

**OpenCL Kernel Language**

Comparison with CUDA and OpenACC

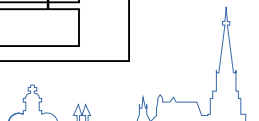
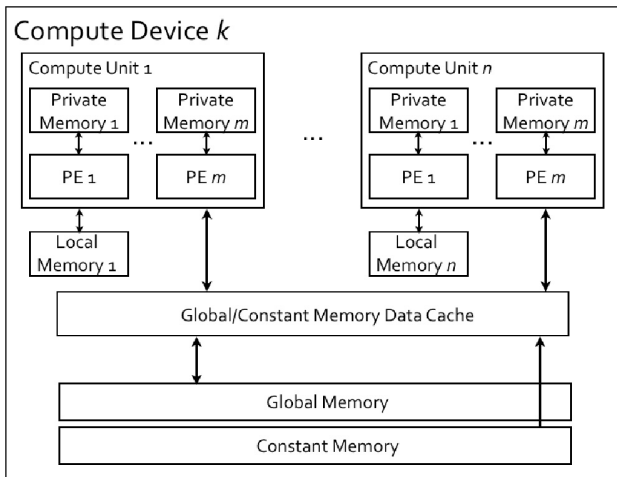
Portable Performance

Summary



# OpenCL Kernel Language

## Primer: OpenCL Device Execution Model

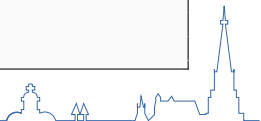


## Sample Operation: Inplace Vector Addition

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} += \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

## OpenCL Kernel

```
__kernel void inplace_add(  
    __global const float * vec1,  
    __global const float * vec2,  
    unsigned int size)  
{  
    for (unsigned int i = get_global_id(0);  
         i < size;  
         i += get_global_size(0))  
        x[i] += y[i];  
}
```



## Just Like C

Datatypes: `float`, `double`, `long`, ... , `float2`, `float4`, ...

Keywords: `for`, `if`, `switch`, ...

...

## Thread Management

ID: `get_local_id(dim)`, `get_global_id(dim)`

Count: `get_local_size(dim)`, `get_global_size(dim)`

Sync: `barrier(flags)`, `mem_fence(flags)`

## Memory Qualifiers

`__global`: Global memory (e.g. GPU RAM)

`__constant`: Constant global memory (vendor-specific!)

`__local`: Shared memory (per workgroup)



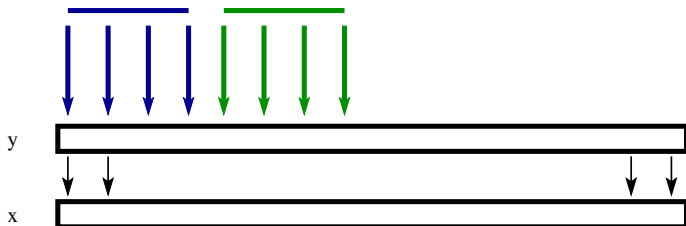
## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$



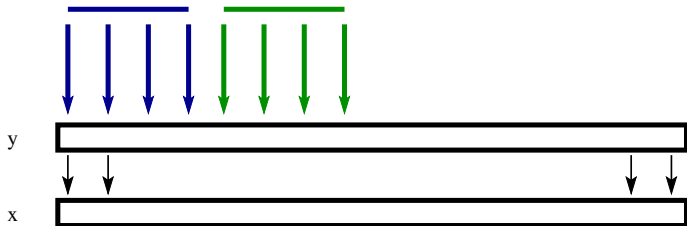
## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$



## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$



## Parameters ( $\gg$ 1000 variations possible)

Local work size, global work size

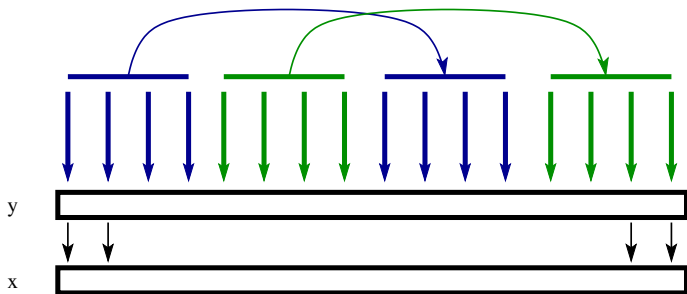
Vector types (float1, float2, ... , float16)

Thread increment type



## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$



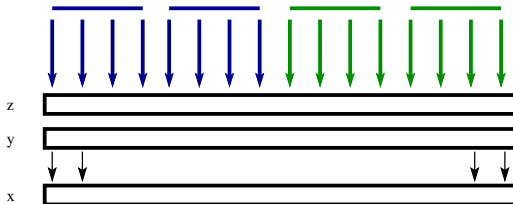
Parameters ( $\gg$  1000 variations possible)

```
for (size_t i = get_global_id(0); i < N;  
     i += get_global_size(0))  
    x[i] = y[i];
```

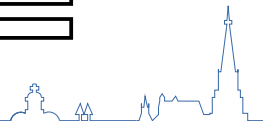
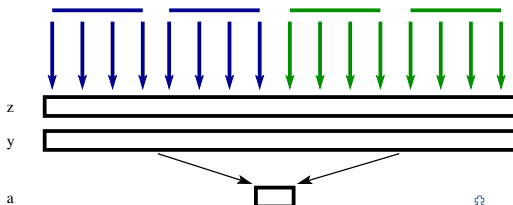




## No Synchronization: Vector Addition



## Synchronization: Dot Product



# Table of Contents

100x Speedup!?

About OpenCL

OpenCL Execution Model

OpenCL Kernel Language

**Comparison with CUDA and OpenACC**

Portable Performance

Summary



## NVIDIA CUDA

```
// GPU kernel:
__global__ void kernel(double *buffer)
{
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    buffer[idx] = 42.0;
}

// host code:
int main()
{
    ...
    cudaMalloc((void**) &buffer, size);
    kernel<<<blocknum, blockDim>>>(buffer);
    ...
}
```

Almost no additional code required

Vendor-lock

Relies on `nvcc` being available (plus version conflicts...)



## OpenCL

```
const char *kernel_string =
"__kernel void mykernel(__global double *buffer) {
    buffer[get_global_id(0)] = 42.0;
};";

int main() {
    ...
    cl_program my_prog = clCreateProgramWithSource(
        my_context, 1, &kernel_string, &source_len, &err);
    cl_build_program(my_prog, 0, NULL, NULL, NULL, NULL);
    cl_kernel my_kernel = clCreateKernel(my_prog,
        "mykernel", &err);
    cl_set_kernel_arg(my_kernel, 0, sizeof(cl_mem), &buffer);
    cl_enqueue_ndrange_kernel(queue, my_kernel, 1, NULL,
        &global_size, &local_size, 0, NULL, NULL);
}
```

Additional boilerplate code required (low-level API)

Broad hardware support (separate SDKs)

Second-best programming model for each vendor



## OpenACC

```
void func(...) {
    #pragma acc data pcopyin(A[0:size][0:size])
    {
        #pragma acc kernels loop
        for(int i=0; i < size; i++)
            for(int j=0; j < size; j++)
                A[i][j] = 42;
    }
}

int main()
{
    double A[1337][1337];
    func(A);
}
```

Simple OpenMP-type pragma annotations

Compiler support?

Insufficient control over memory transfers?



## Thread Explosion Problem

```
void func(double *A) {
    #pragma omp parallel for
    for(int i=0; i<1337; i++) A[i] = 42.0;
}

void func2(double *A) {
    #pragma omp parallel for
    for(int i=0; i<1337; i++) func(A);
}

int main() {
    double A[1337];
    func(A); // okay
    func2(A); // boom...
}
```

Usually not as obvious as here  
Problem when OpenMP'ing MPI code  
Headache for library implementors



# Table of Contents

100x Speedup!?

About OpenCL

OpenCL Execution Model

OpenCL Kernel Language

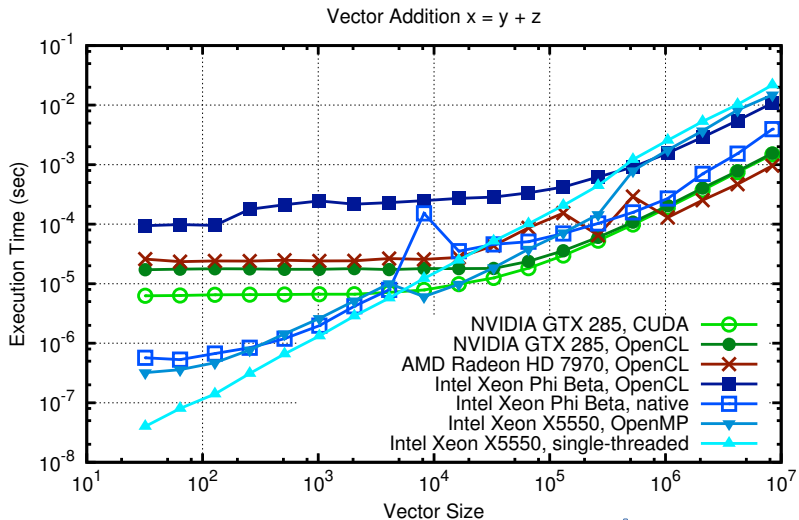
Comparison with CUDA and OpenACC

**Portable Performance**

Summary

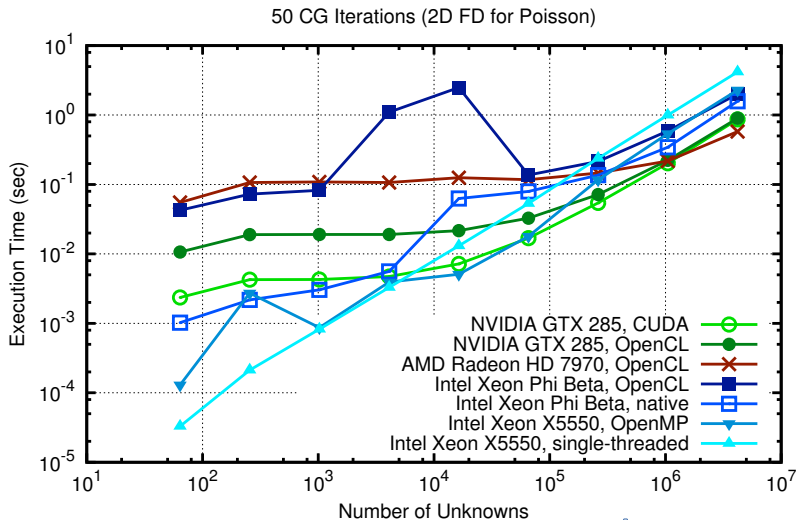


# Benchmarks





# Benchmarks



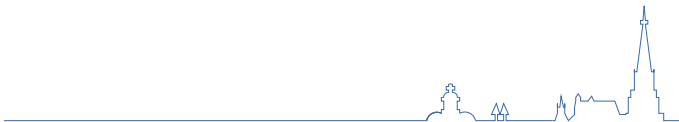
## Scope for Portability Study

Vector and matrix-vector operations (BLAS levels 1 and 2)

Limited by memory bandwidth

## Key Question (Memory-Bandwidth-Limited Kernels)

Good performance of complicated kernels  
by optimizing the simplest kernel?



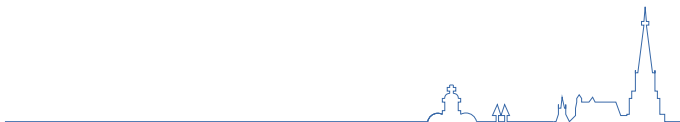
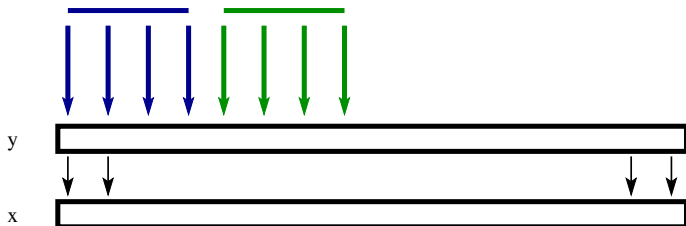
## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$



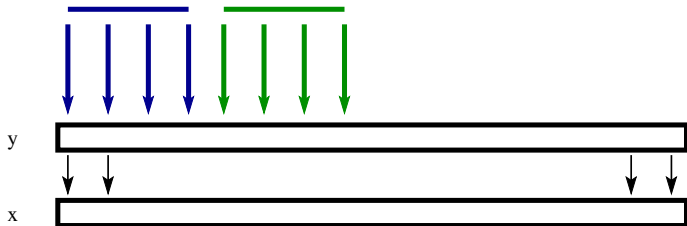
## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$



## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$



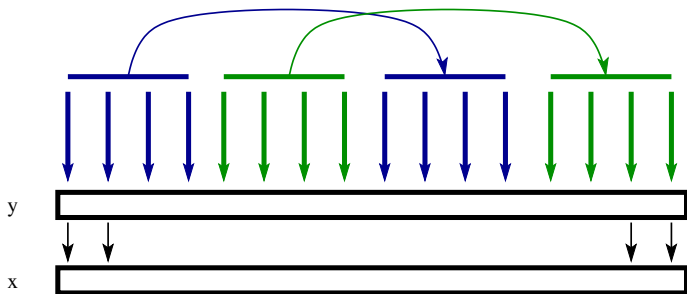
## Parameters (1900 variations)

- Local work size, global work size
- Vector types (float1, float2, ... , float16)
- Thread increment type



## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$



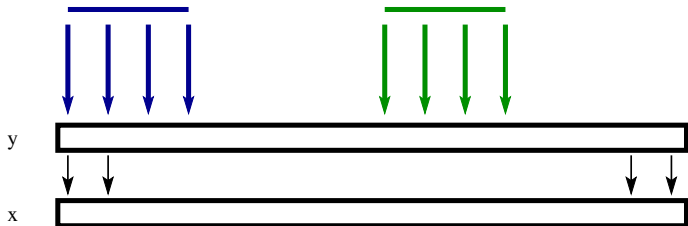
## Parameters (1900 variations)

```
for (size_t i = get_global_id(0); i < N;  
     i += get_global_size(0))  
    x[i] = y[i];
```



## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$

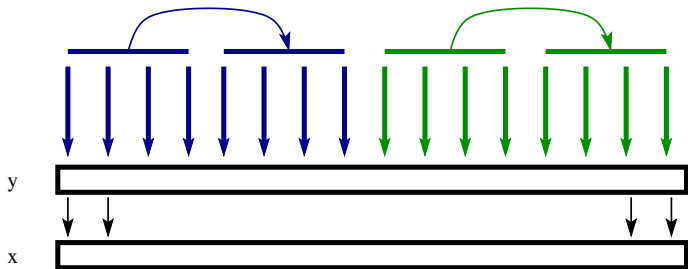


## Parameters (1900 variations)

```
for (size_t i = group_start + get_local_id(0);  
     i < group_end;  
     i += get_local_size(0)) x[i] = y[i];
```

## Vector Assignment (Copy) Kernel

$x \leftarrow y$  for (large) vectors  $x, y$



## Parameters (1900 variations)

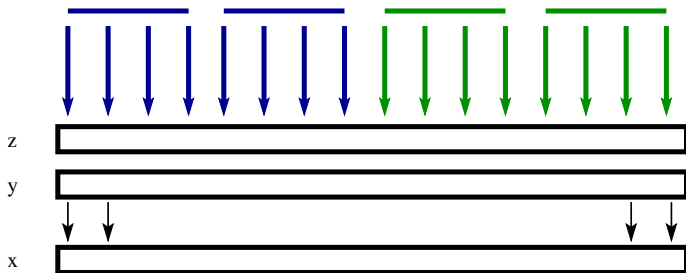
```
for (size_t i = group_start + get_local_id(0);  
    i < group_end; i += get_local_size(0))  
    x[i] = y[i];
```



## Operations

Vector copy, vector addition, inner product

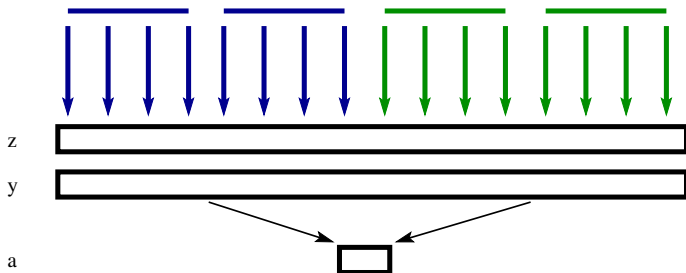
Matrix-vector product



## Operations

Vector copy, vector addition, inner product

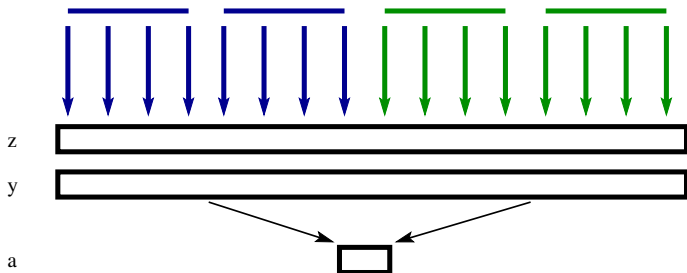
Matrix-vector product



## Operations

Vector copy, vector addition, inner product

Matrix-vector product



## Devices

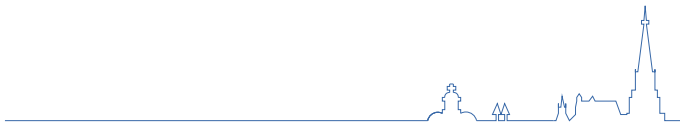
AMD: Radeon HD 5850, FirePro W9000

INTEL: Dual Socket Xeon E5-2670, Xeon Phi

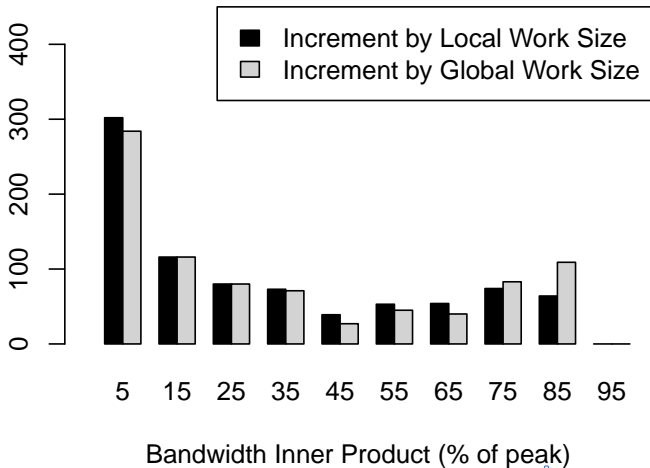
NVIDIA: GTX 285, Tesla K20m



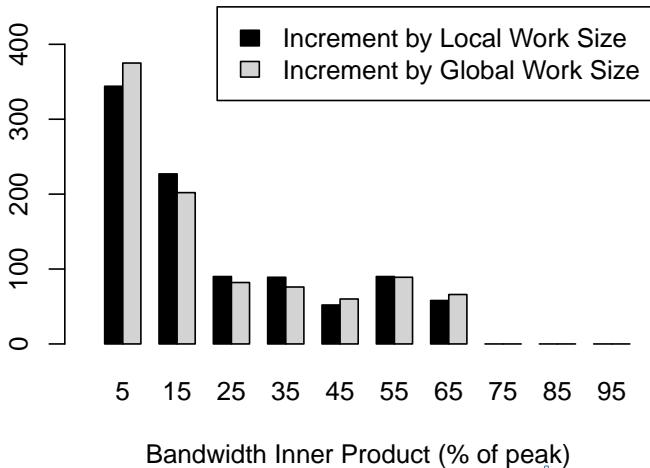
## Histograms



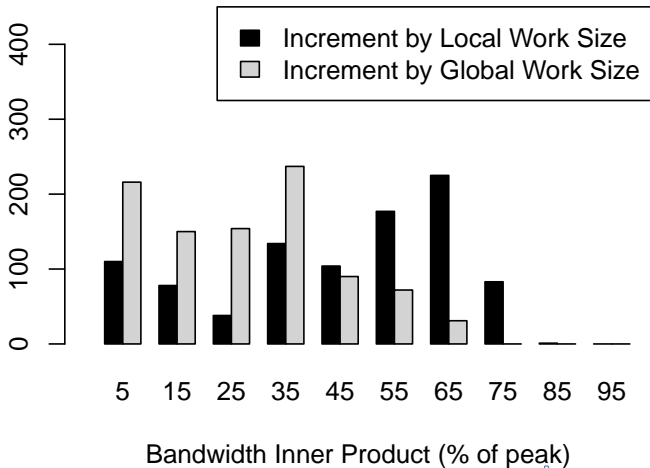
## AMD FirePro W9000



## NVIDIA Tesla K20m

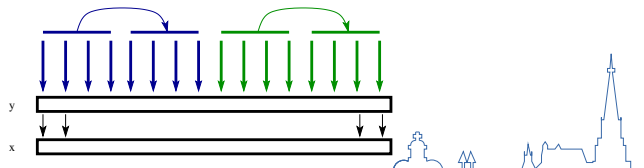
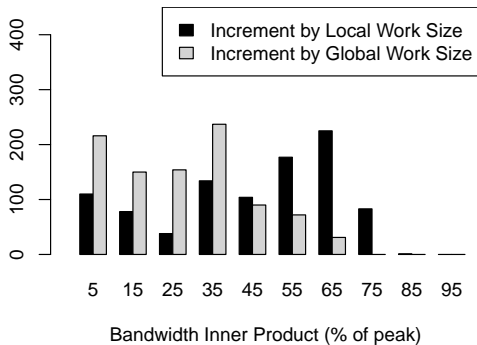


## 2x INTEL Xeon E5-2670



# Portable Performance

## 2x INTEL Xeon E5-2670



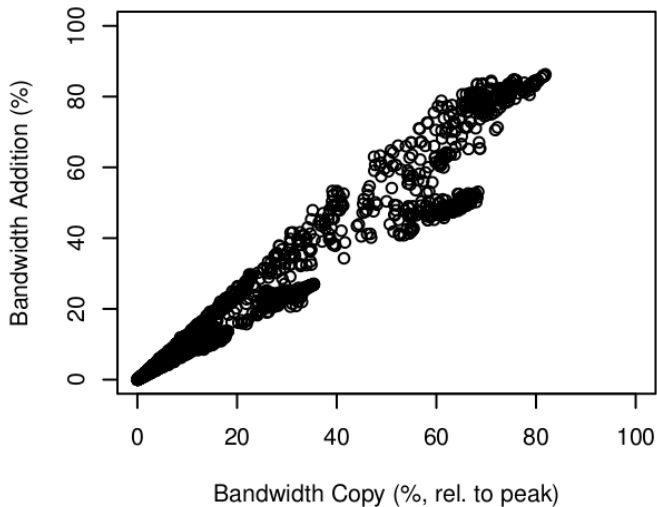


**[Addition|Inner Product|Matrix-Vector] vs. Copy Kernel**

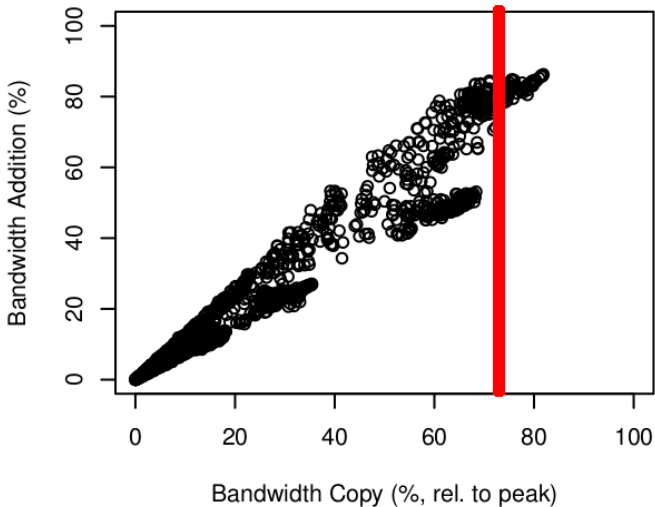
Same Device



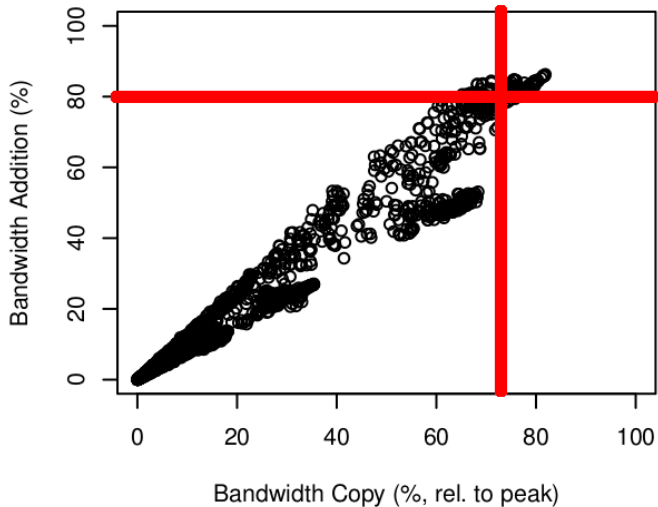
## NVIDIA GeForce GTX 285



## NVIDIA GeForce GTX 285

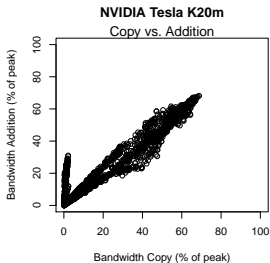


## NVIDIA GeForce GTX 285

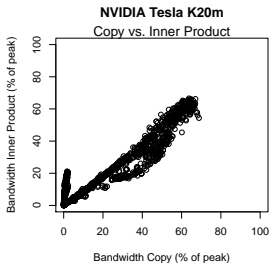


## NVIDIA Tesla K20m

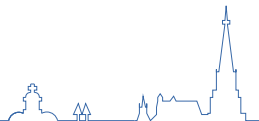
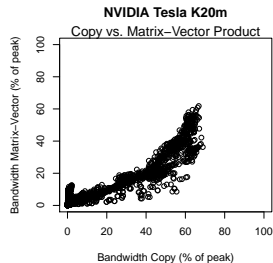
### Addition



### Inner Product

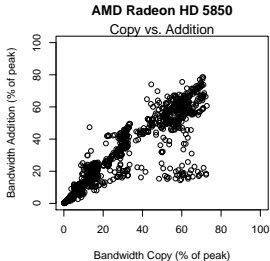


### Mat-Vec Product

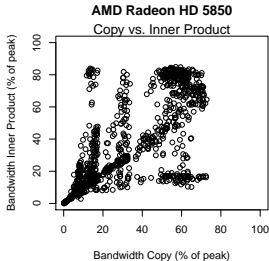


## AMD Radeon HD 5850

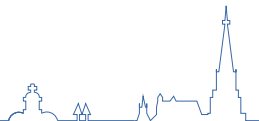
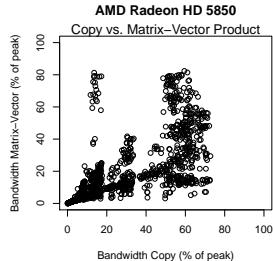
### Addition



### Inner Product

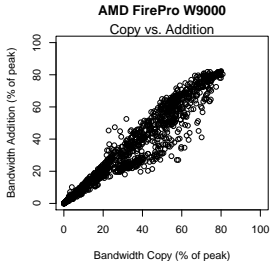


### Mat-Vec Product

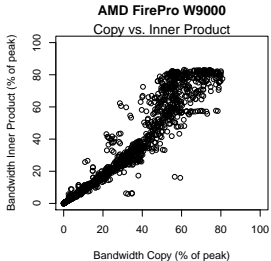


## AMD FirePro W9000

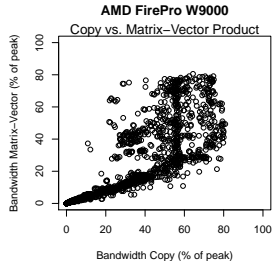
### Addition



### Inner Product

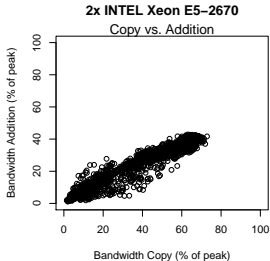


### Mat-Vec Product

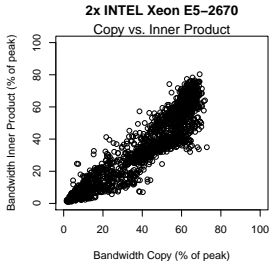


## INTEL Dual Xeon E5-2670

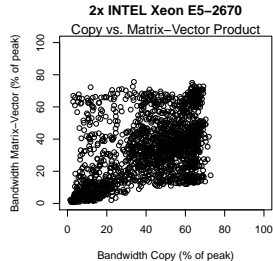
### Addition



### Inner Product



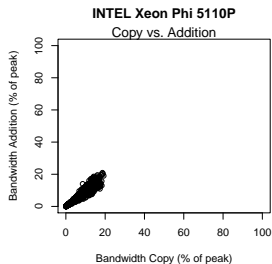
### Mat-Vec Product



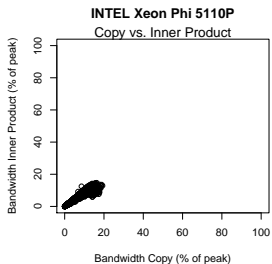


## INTEL Xeon Phi

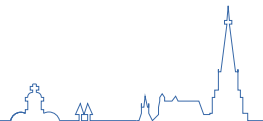
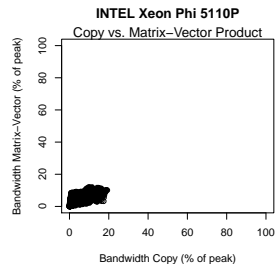
### Addition



### Inner Product



### Mat-Vec Product

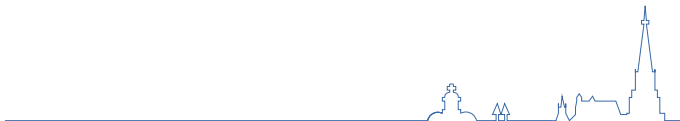


Conclusio:

Focus on fastest configurations for copy-kernel sufficient

Good choice:

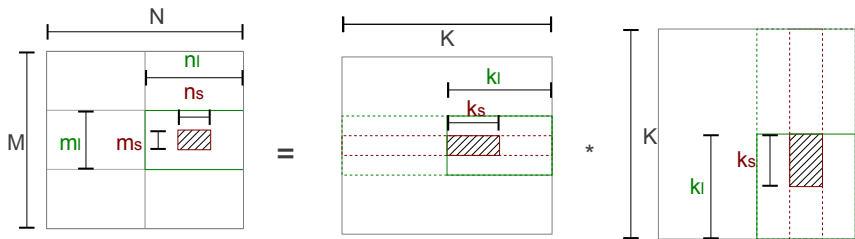
Local workgroup size of 128 with 128 workgroups

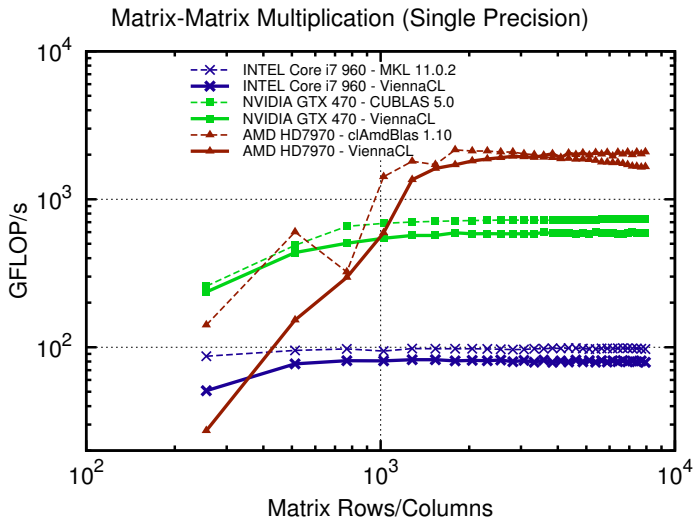


## Matrix-Matrix Multiplication

Compute-bound

Block-decomposition to maximize cache utilization





## OpenCL

- Program CPUs and accelerators (GPUs, MIC, etc.) from different vendors
- Clean integration into existing projects
- Not a silver bullet

## Recommendations

- Program with data locality and data movement in mind
- Compilers cannot fix bad programming
- There is no free lunch

## Convenient Use of OpenCL

- Use libraries built on top of OpenCL

Tue, July 15, 9:00am:

*Lessons Learned in Developing the Linear Algebra Library ViennaCL*

