

Performance Portability Study of Linear Algebra Kernels in OpenCL

Karl Rupp^{1,2}, Philippe Tillet¹, Florian Rudolf¹,
Josef Weinbub¹, Ansgar Jüngel², Tibor Grasser¹

`rupp@iue.tuwien.ac.at`

 @karlrupp



¹ Institute for Microelectronics, TU Wien, Austria

² Institute for Analysis and Scientific Computing, TU Wien, Austria

International Workshop on OpenCL 2014

May 13th, 2014



OpenCL Code Library Expectations



OpenCL Code Library Expectations

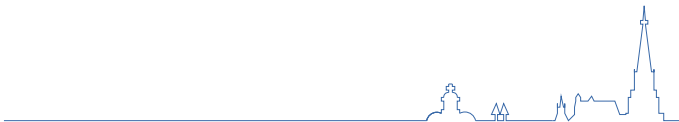
Just work on my (user) hardware

Run fast on my (user) hardware

Integrate easily into my (user) code

Mix with CUDA and OpenMP

Available for free



OpenCL Code Library Expectations

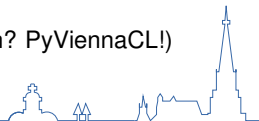
- Just work on my (user) hardware
- Run fast on my (user) hardware
- Integrate easily into my (user) code
- Mix with CUDA and OpenMP
- Available for free

Work Based on Developer Experiences



<http://viennacl.sourceforge.net/>

Free open-source C++ linear algebra library (Python? PyViennaCL!)
API similar to Boost.uBLAS



How OpenCL Helps

Unified hardware information queries

Unified kernel language



How OpenCL Helps

- Unified hardware information queries

- Unified kernel language

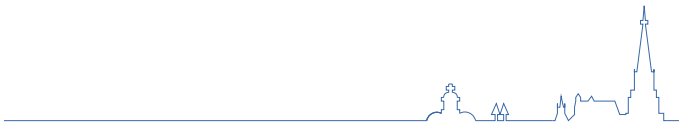
Developer Tasks

- OpenCL \neq automatic performance portability

- Parameterize each kernel

- Performance models vs. (auto)tuning

- Impractical to repeat for each kernel



Scope for Portability Study

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

Limited by memory bandwidth

Matrix-matrix-multiplication (only briefly today)



Scope for Portability Study

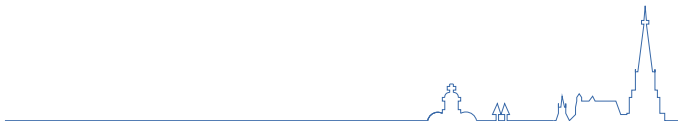
Vector and matrix-vector operations (BLAS levels 1 and 2)

Limited by memory bandwidth

Matrix-matrix-multiplication (only briefly today)

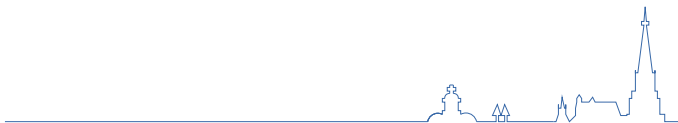
Key Question (Memory-Bandwidth-Limited Kernels)

Good performance of complicated kernels
by optimizing the simplest kernel?



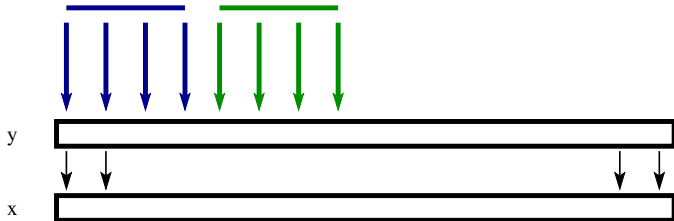
Vector Assignment (Copy) Kernel

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



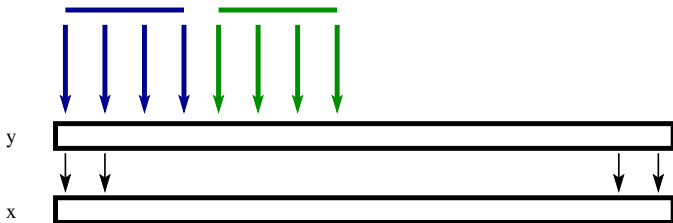
Vector Assignment (Copy) Kernel

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



Vector Assignment (Copy) Kernel

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

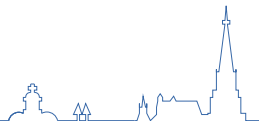


Parameters (1900 variations)

Local work size, global work size

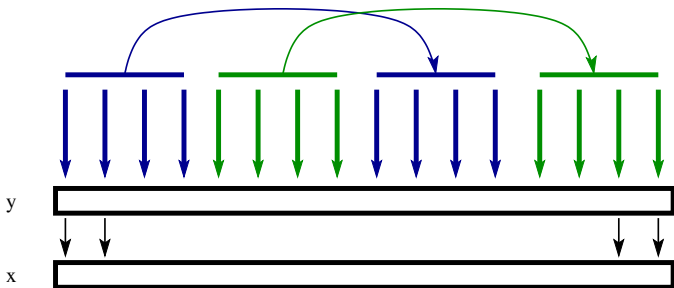
Vector types (float1, float2, ... , float16)

Thread increment type



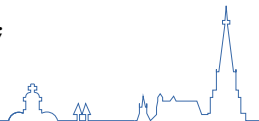
Vector Assignment (Copy) Kernel

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



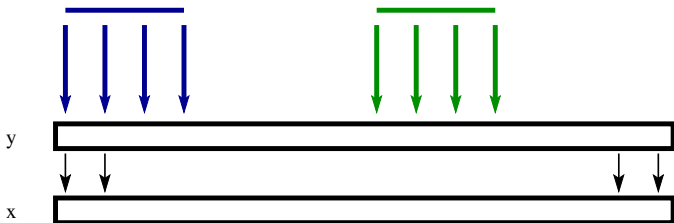
Parameters (1900 variations)

```
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

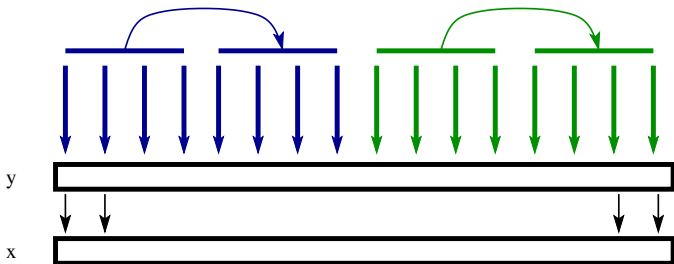


Parameters (1900 variations)

```
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



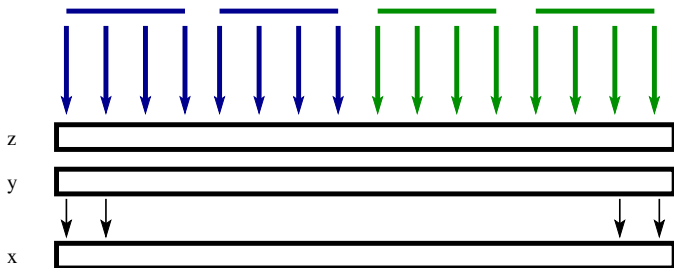
Parameters (1900 variations)

```
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

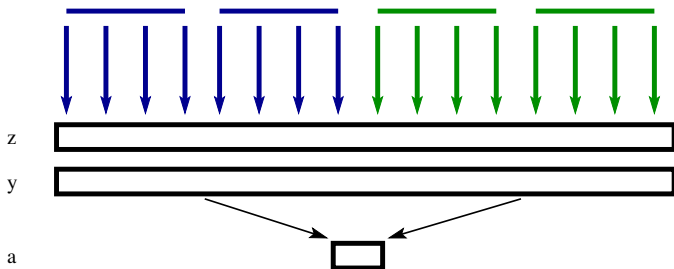
Matrix-vector product



Operations

Vector copy, vector addition, inner product

Matrix-vector product

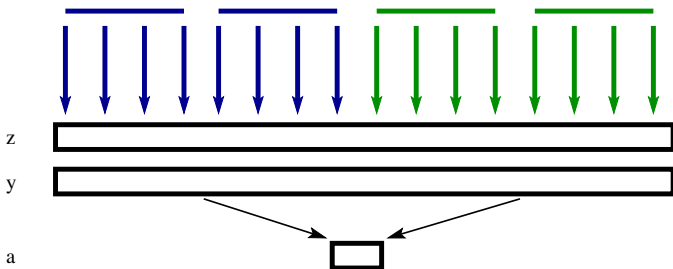


Benchmark Setting

Operations

Vector copy, vector addition, inner product

Matrix-vector product

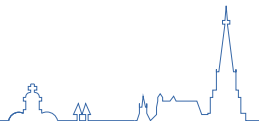


Devices

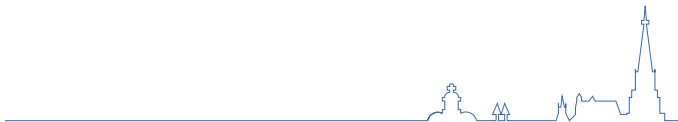
AMD: A10-5800 APU, HD 5850 GPU

INTEL: Dual Socket Xeon E5-2670, Xeon Phi

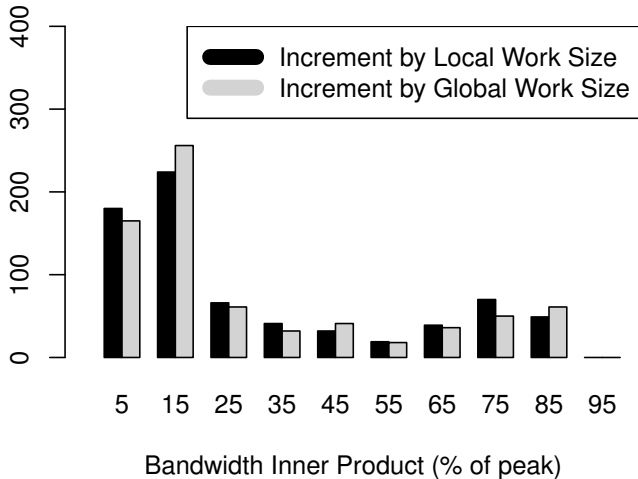
NVIDIA: GTX 285, Tesla K20m



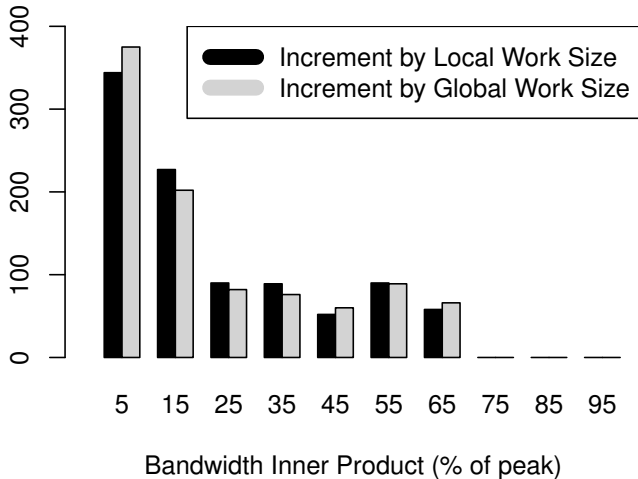
Histograms



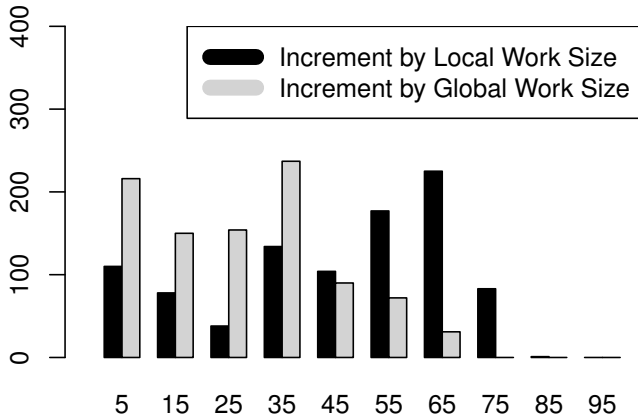
AMD Radeon HD 5850



NVIDIA Tesla K20m



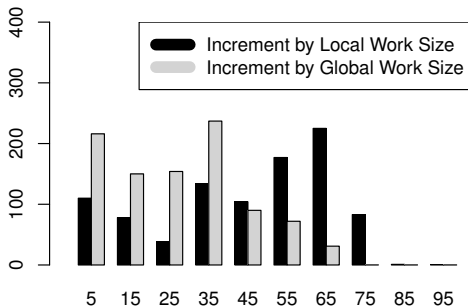
Intel Xeon E5-2670



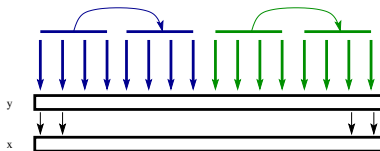
Bandwidth Inner Product (% of theoretical peak)

Benchmark

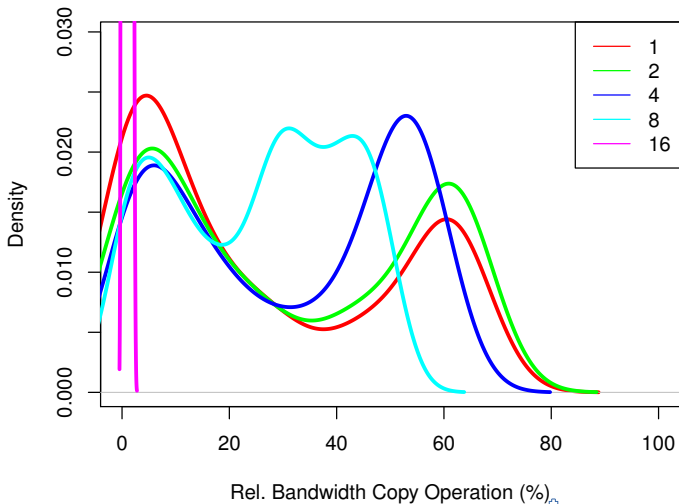
Intel Xeon E5-2670



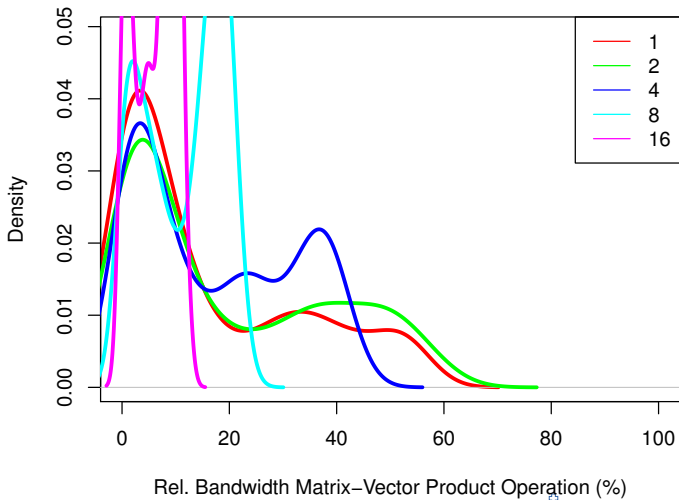
Bandwidth Inner Product (% of theoretical peak)



NVIDIA Tesla K20m



NVIDIA Tesla K20m

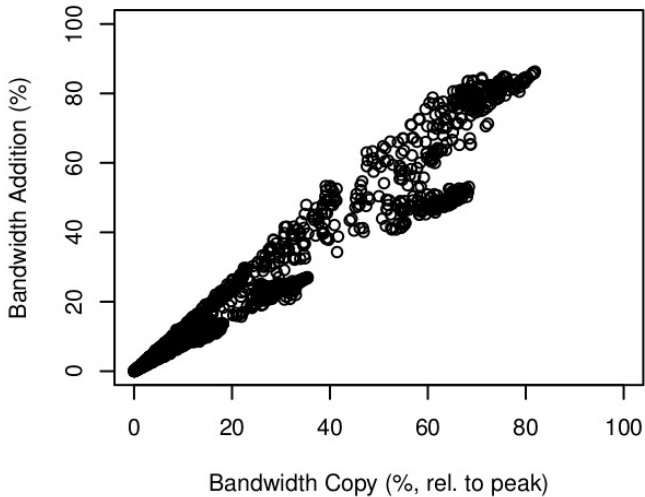


[Addition|Inner Product|Matrix-Vector] vs. Copy Kernel

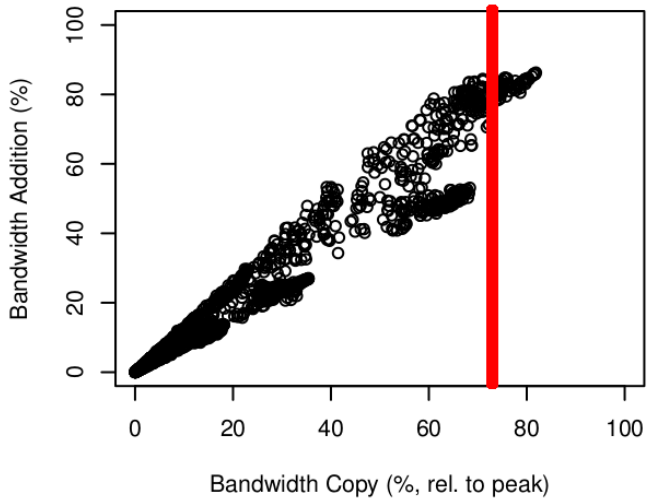
Same Device



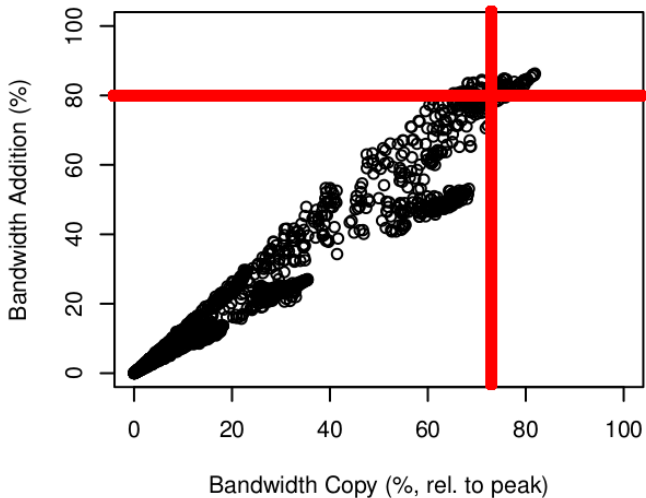
NVIDIA GeForce GTX 285



NVIDIA GeForce GTX 285

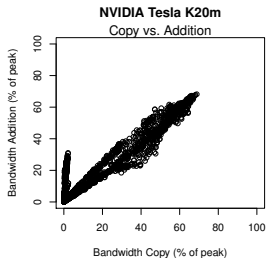


NVIDIA GeForce GTX 285

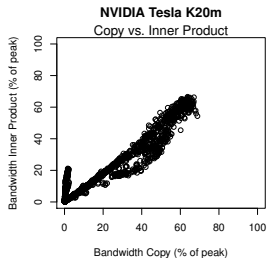


NVIDIA Tesla K20m

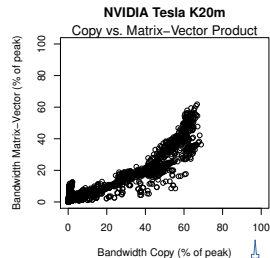
Addition



Inner Product

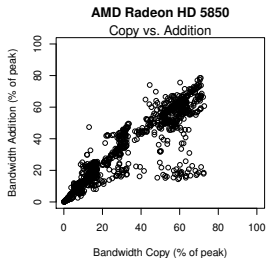


Mat-Vec Product

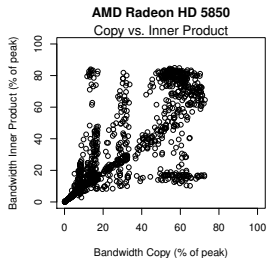


AMD Radeon HD 5850

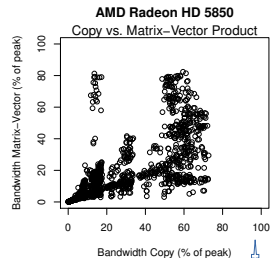
Addition



Inner Product

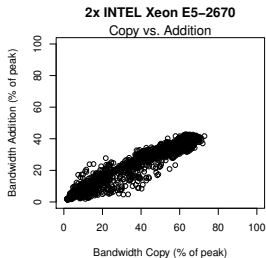


Mat-Vec Product

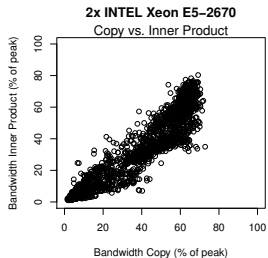


INTEL Dual Xeon E5-2670

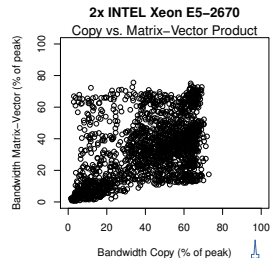
Addition



Inner Product

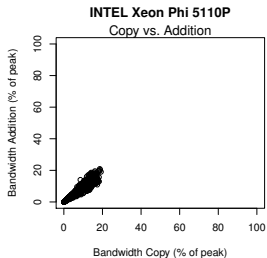


Mat-Vec Product

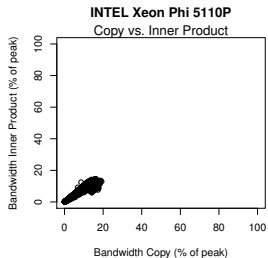


INTEL Xeon Phi

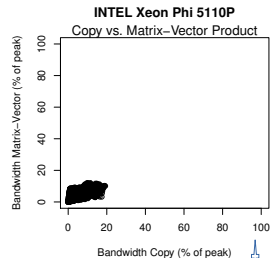
Addition



Inner Product



Mat-Vec Product



Conclusio:

Focus on fastest configurations for copy-kernel sufficient

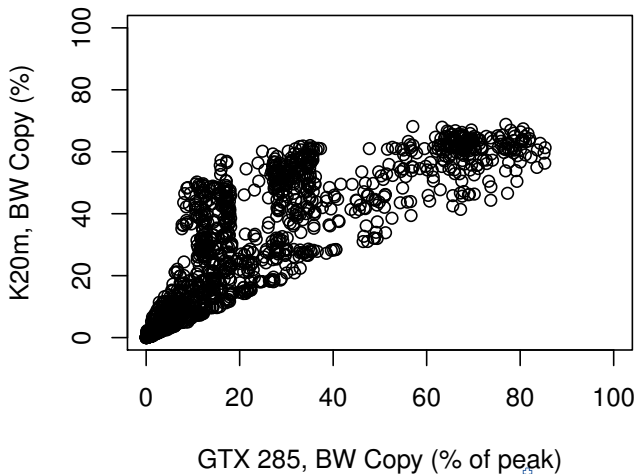


[Copy|Addition|Inner Product|Matrix-Vector] vs. Copy Kernel

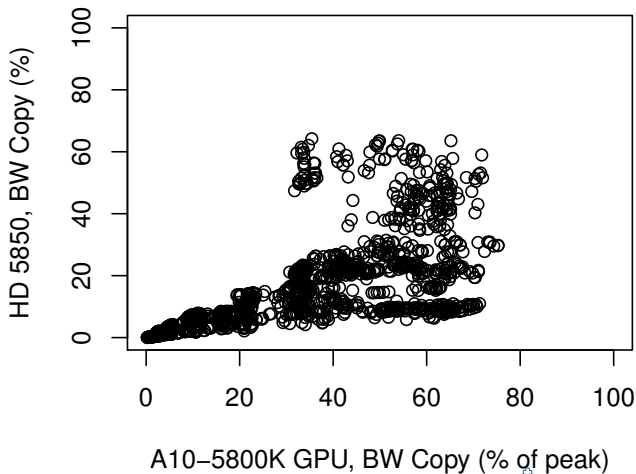
Different Device, Same Vendor



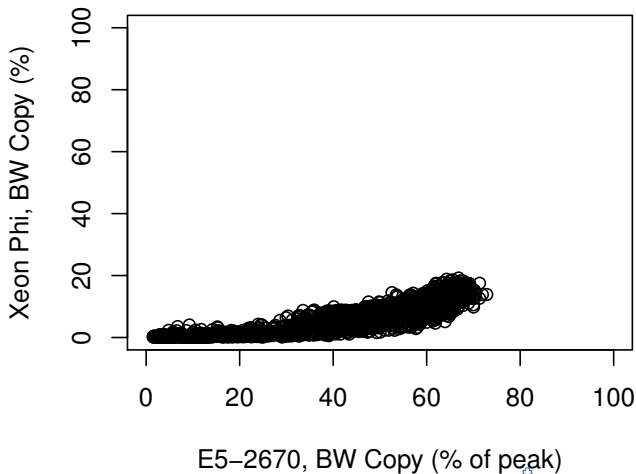
NVIDIA Hardware (x: GTX 285, y: K20m)



AMD Hardware (x: A10-5800K GPU, y: HD 5850)

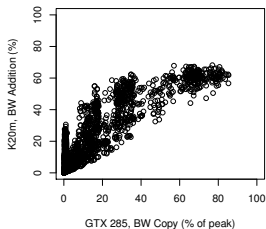


INTEL Hardware (x: Xeon Phi, E5-2670)

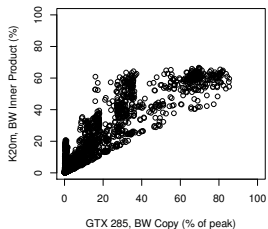


NVIDIA Hardware (x: GTX 285, y: K20m)

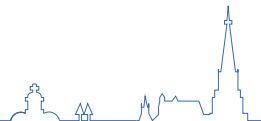
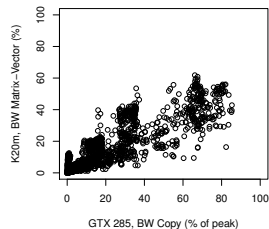
Addition



Inner Product

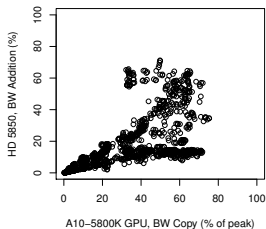


Mat-Vec Product

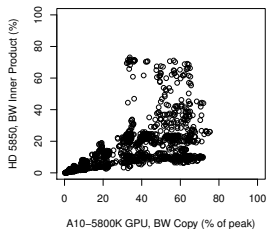


AMD Hardware (x: A10-5800K GPU, y: HD 5850)

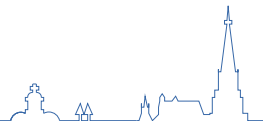
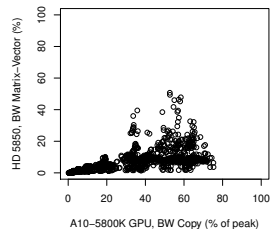
Addition



Inner Product



Mat-Vec Product



Conclusio:

Certain Performance Portability per Vendor



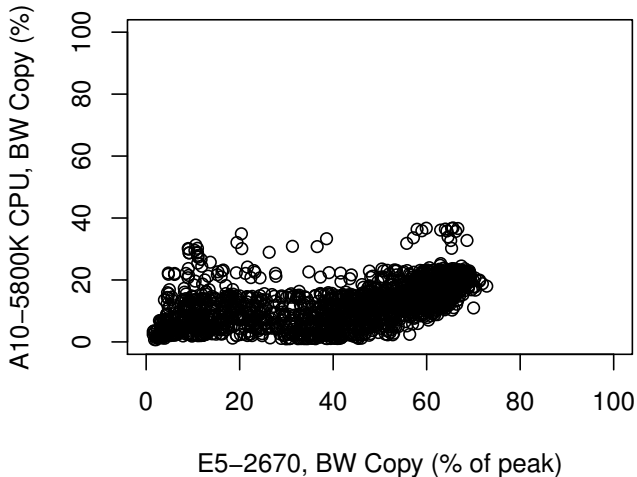
[Copy|Addition|Inner Product|Matrix-Vector] vs. Copy Kernel

Different Device, Different Vendor

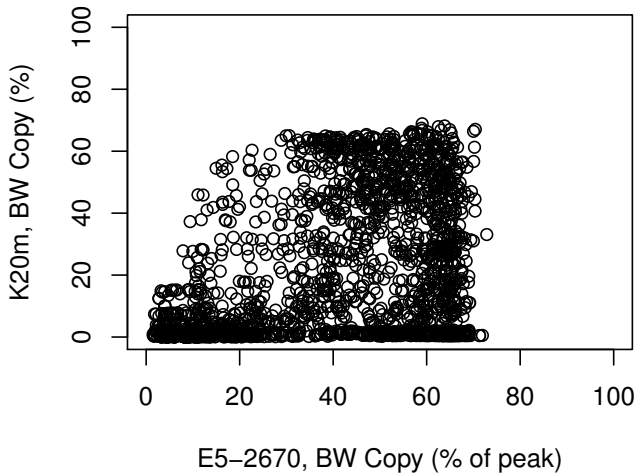


Benchmark

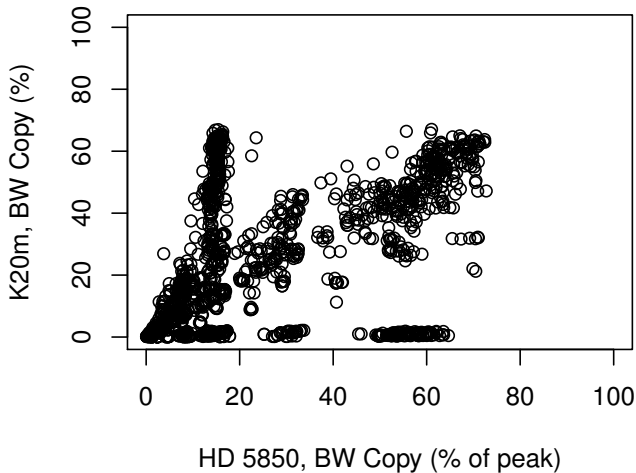
x: INTEL CPU, y: AMD CPU



x: INTEL CPU, y: NVIDIA GPU

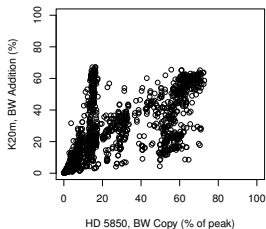


x: AMD GPU, y: NVIDIA GPU

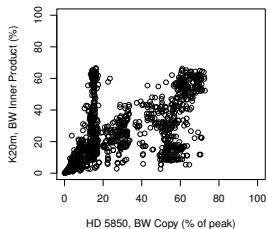


x: AMD HD 5850, y: NVIDIA K20m

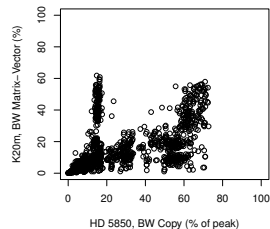
Addition



Inner Product



Mat-Vec Product



Conclusio:

Fast Configurations Across Vendors Exist

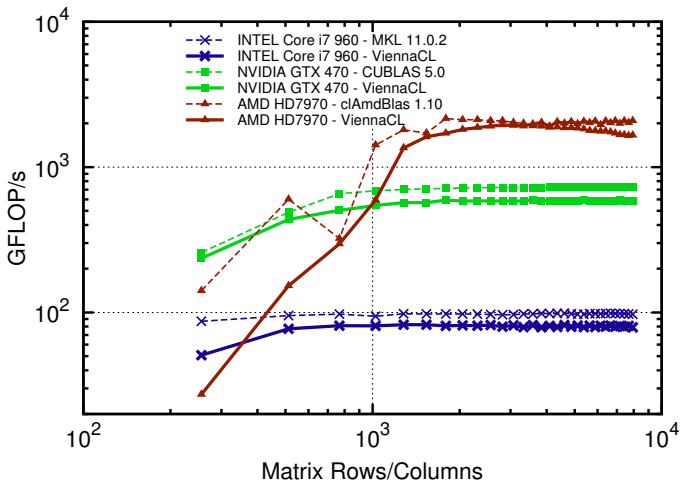


Matrix-Matrix Multiplication

Single Precision

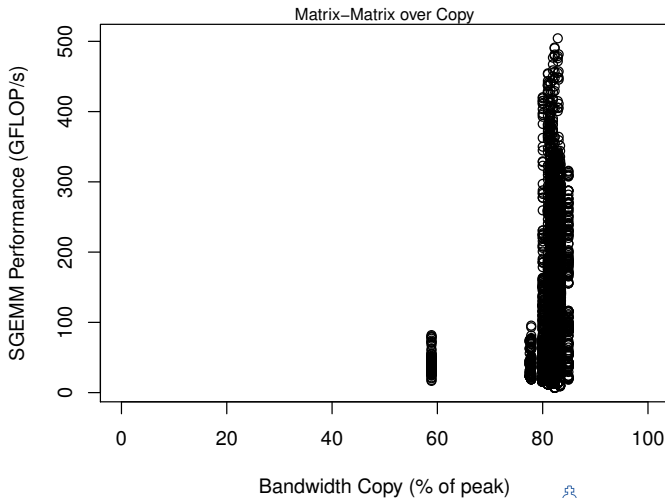


Benchmark



Ph. Tillet *et al.*: HotPar'13

NVIDIA GeForce 470 GTX



Copy Kernel

Prototype for complicated memory-bandwidth-limited kernels

Strong performance correlation

Reduces search space for auto-tuning



Copy Kernel

Prototype for complicated memory-bandwidth-limited kernels

Strong performance correlation

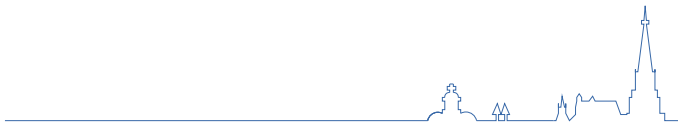
Reduces search space for auto-tuning

Implications

Tune primarily with memory transfers in mind

Prefer regular memory access patterns

Use *appropriate* vector data types



Copy Kernel

- Prototype for complicated memory-bandwidth-limited kernels
- Strong performance correlation
- Reduces search space for auto-tuning

Implications

- Tune primarily with memory transfers in mind
- Prefer regular memory access patterns
- Use *appropriate* vector data types

Next Steps

- Build a device parameter database
- Provide tools for crowd-tuning

