

Writing Performance-Portable Code for GPUs

Karl Rupp^{1,2}

`rupp@iue.tuwien.ac.at`

 @karlrupp

with contributions from

Philippe Tillet¹, Florian Rudolf¹,
Josef Weinbub¹, Ansgar Jüngel², Tibor Grasser¹
(based on stimuli from PETSc+ViennaCL users)



¹ Institute for Microelectronics, TU Wien, Austria

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

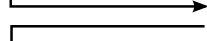
Introduction



Conjugate Gradient - Standard Formulation



Conjugate Gradient - Kernel Parameters



Parameter Study for Portable Performance



Conjugate Gradient - Pipelining



BiCGStab and GMRES - Benchmark Results



Conclusion



Positions

PhD student at TU Wien (2009-2011)

Postdoc at ANL (09/2012-09/2013)

Postdoc at TU Wien (01/2012-09/2012, 09/2013-current)

Research Interests

Semiconductor device simulation

Numerical solution of PDEs

Parallel computing

Software Development

PETSc

ViennaCL

ViennaSHE

...



Iterative Solvers

Matrix-vector products and vector operations only

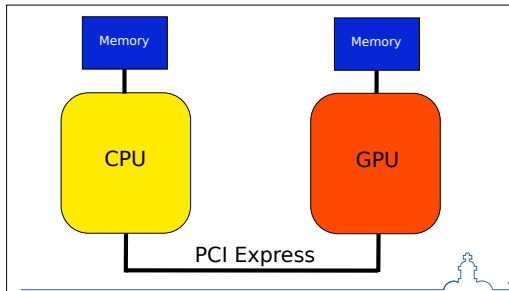
Expose more fine-grained parallelism

Preconditioners often desirable

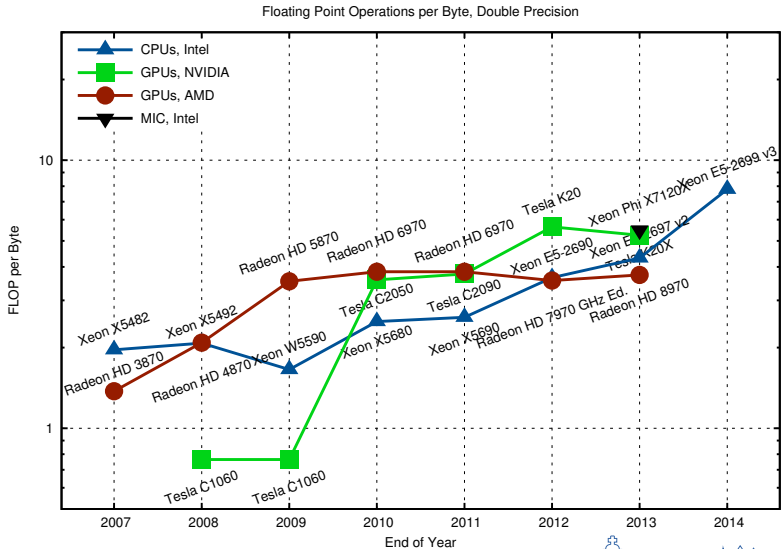
Accelerators (CUDA, OpenCL)

Graphics processing units (GPUs)

Intel Xeon Phi



Introduction



Pseudocode

Choose x_0

$$p_0 = r_0 = b - Ax_0$$

For $i = 0$ until convergence

1. Compute and store Ap_i
2. Compute $\langle p_i, Ap_i \rangle$
3. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
4. $x_{i+1} = x_i + \alpha_i p_i$
5. $r_{i+1} = r_i - \alpha_i Ap_i$
6. Compute $\langle r_{i+1}, r_{i+1} \rangle$
7. $\beta_i = \langle r_{i+1}, r_{i+1} \rangle / \langle r_i, r_i \rangle$
8. $p_{i+1} = r_{i+1} + \beta_i p_i$

EndFor

BLAS-based Implementation

-

SpMV, AXPY

For $i = 0$ until convergence

1. SpMV
2. DOT
3. -
4. AXPY
5. AXPY
6. DOT
7. -
8. AXPY

EndFor



Pseudocode

Choose x_0

$$p_0 = r_0 = b - Ax_0$$

For $i = 0$ until convergence

1. Compute and store Ap_i
2. Compute $\langle p_i, Ap_i \rangle$
3. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
4. $x_{i+1} = x_i + \alpha_i p_i$
5. $r_{i+1} = r_i - \alpha_i Ap_i$
6. Compute $\langle r_{i+1}, r_{i+1} \rangle$
7. $\beta_i = \langle r_{i+1}, r_{i+1} \rangle / \langle r_i, r_i \rangle$
8. $p_{i+1} = r_{i+1} + \beta_i p_i$

EndFor

BLAS-based Implementation

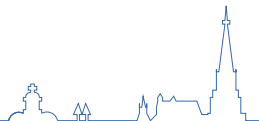
-

SpMV, AXPY

For $i = 0$ until convergence

1. SpMV
2. DOT \leftarrow Global sync!
3. -
4. AXPY
5. AXPY
6. DOT \leftarrow Global sync!
7. -
8. AXPY

EndFor



Pseudocode

Choose x_0

$$p_0 = r_0 = b - Ax_0$$

For $i = 0$ until convergence

1. Compute and store Ap_i
2. Compute $\langle p_i, Ap_i \rangle$
3. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
4. $x_{i+1} = x_i + \alpha_i p_i$
5. $r_{i+1} = r_i - \alpha_i Ap_i$
6. Compute $\langle r_{i+1}, r_{i+1} \rangle$
7. $\beta_i = \langle r_{i+1}, r_{i+1} \rangle / \langle r_i, r_i \rangle$
8. $p_{i+1} = r_{i+1} + \beta_i p_i$

EndFor

BLAS-based Implementation

-

SpMV, AXPY

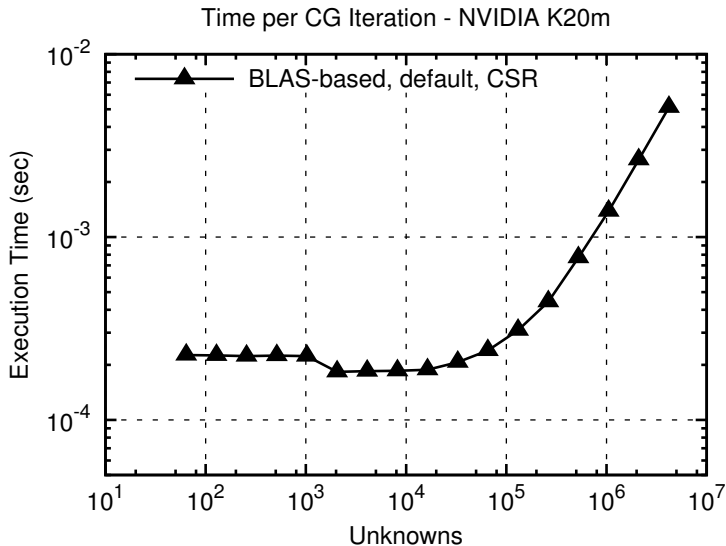
For $i = 0$ until convergence

1. SpMV \leftarrow No caching of Ap_i
2. DOT \leftarrow Global sync!
3. -
4. AXPY
5. AXPY \leftarrow No caching of r_{i+1}
6. DOT \leftarrow Global sync!
7. -
8. AXPY

EndFor



Conjugate Gradients



(2D Finite Difference Discretization)

Implications

Kernel launches expensive

Delicate balance for preconditioners

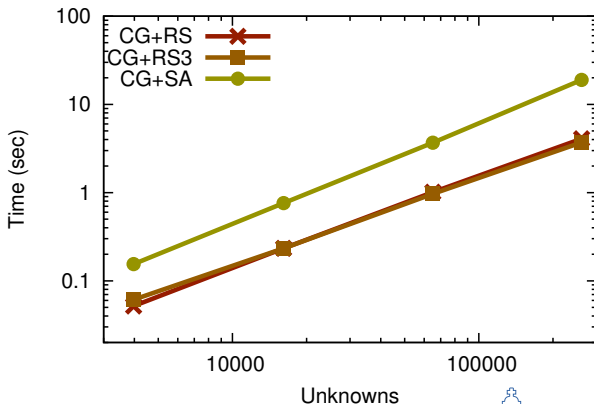


Implications

Kernel launches expensive

Delicate balance for preconditioners

Comparison of Execution Times on CPU

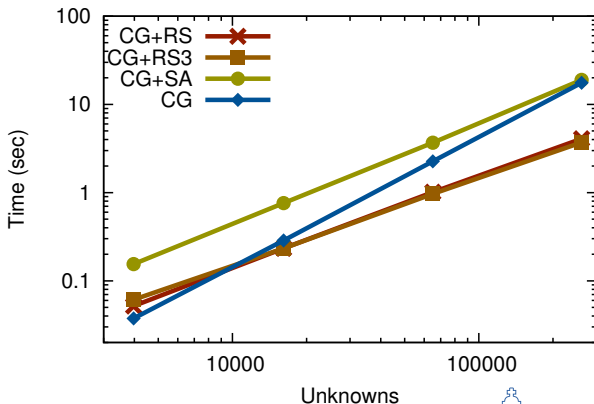


Implications

Kernel launches expensive

Delicate balance for preconditioners

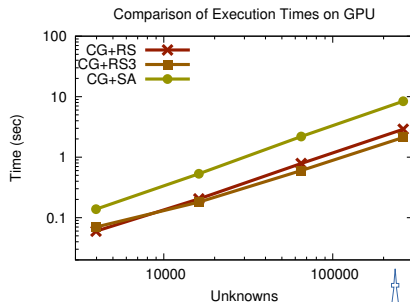
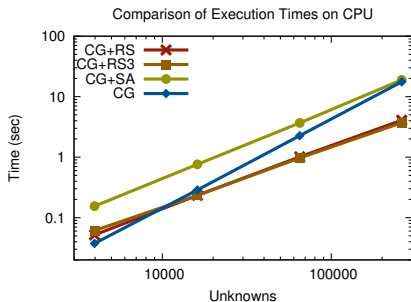
Comparison of Execution Times on CPU



Implications

Kernel launches expensive

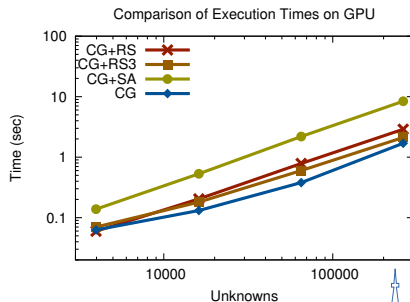
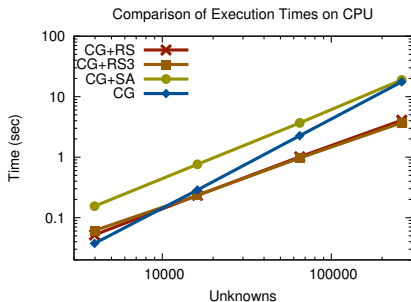
Delicate balance for preconditioners



Implications

Kernel launches expensive

Delicate balance for preconditioners

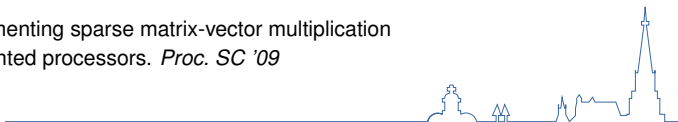


Optimization 1

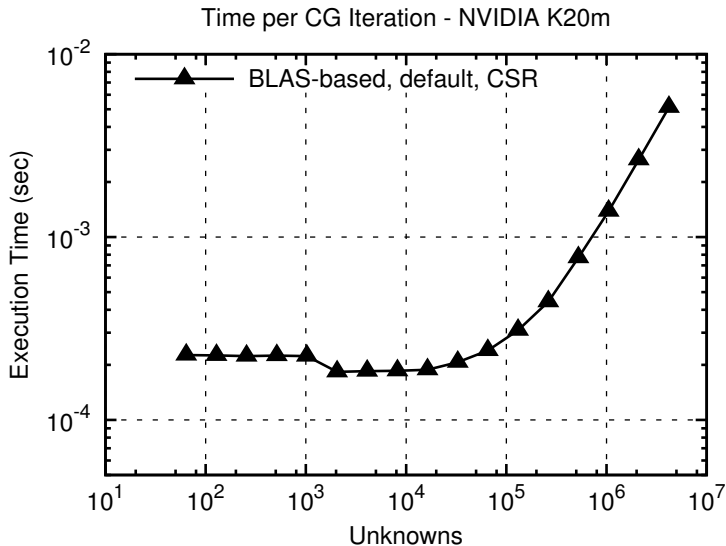
Get best performance out of SpMV

Compare different sparse matrix types

Cf.: N. Bell: Implementing sparse matrix-vector multiplication on throughput-oriented processors. *Proc. SC '09*

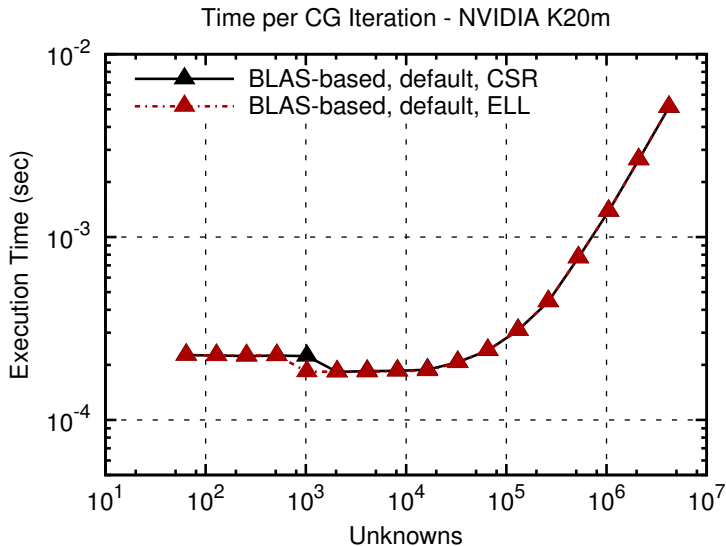


Conjugate Gradients



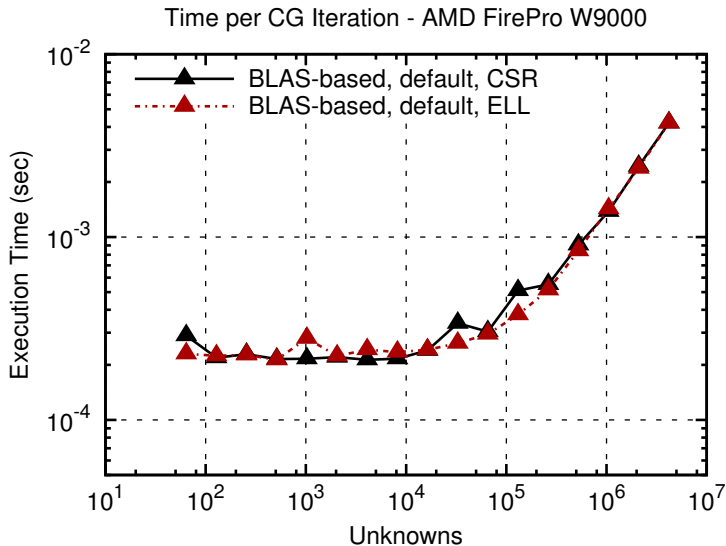
(2D Finite Difference Discretization)

Conjugate Gradients



(2D Finite Difference Discretization)

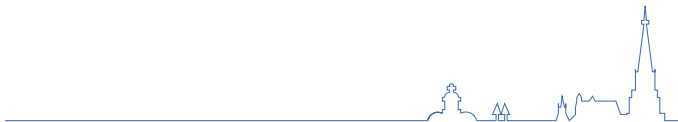
Conjugate Gradients



(2D Finite Difference Discretization)

Optimization 2

Optimize kernel parameters for each operation



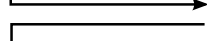
Introduction



Conjugate Gradient - Standard Formulation



Conjugate Gradient - Kernel Parameters



Parameter Study for Portable Performance



Conjugate Gradient - Pipelining



BiCGStab and GMRES - Benchmark Results



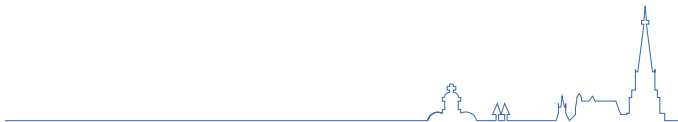
Conclusion



Scope for Portability Study

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

Limited by memory bandwidth



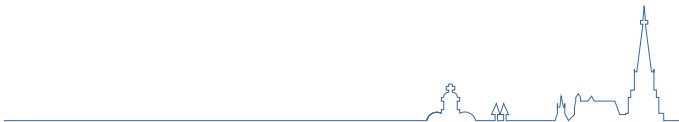
Scope for Portability Study

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?



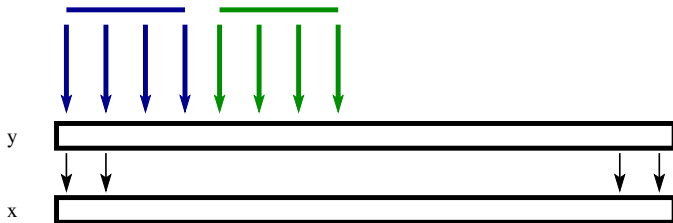
Vector Assignment (Copy) Kernel

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



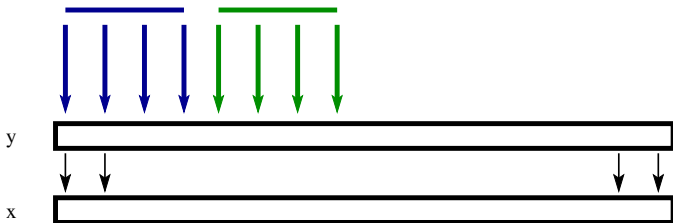
Vector Assignment (Copy) Kernel

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



Vector Assignment (Copy) Kernel

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

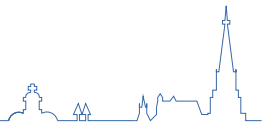


Parameters (1900 variations)

Local work size, global work size

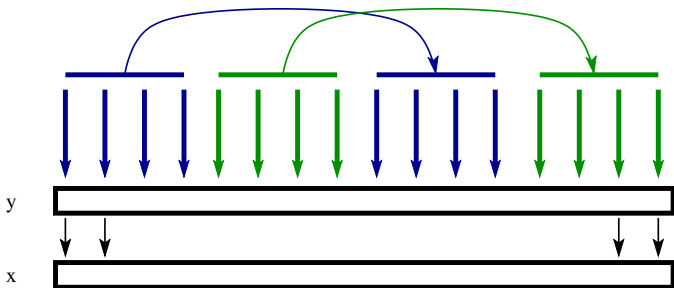
Vector types (float1, float2, ... , float16)

Thread increment type



Vector Assignment (Copy) Kernel

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



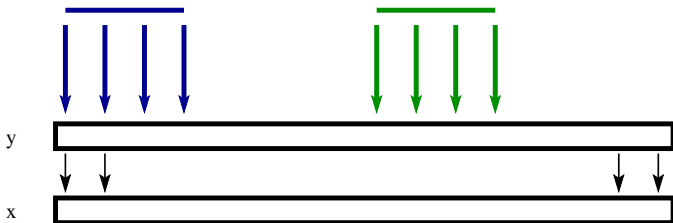
Parameters (1900 variations)

```
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

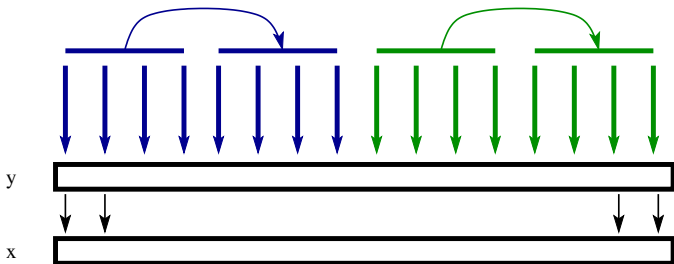


Parameters (1900 variations)

```
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



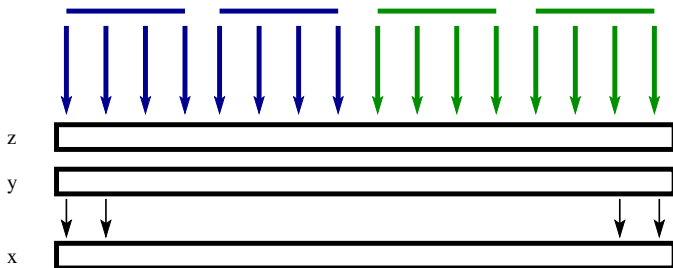
Parameters (1900 variations)

```
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

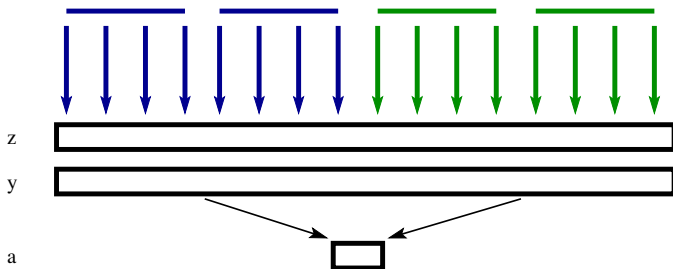
Matrix-vector product



Operations

Vector copy, vector addition, inner product

Matrix-vector product

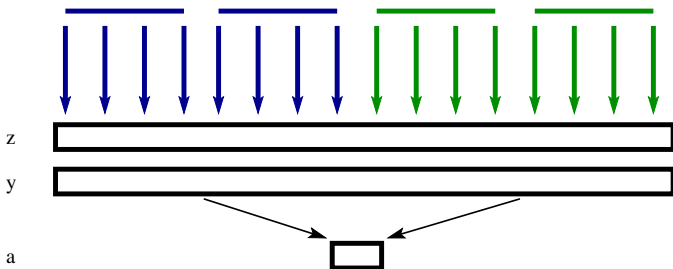


Benchmark Setting

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

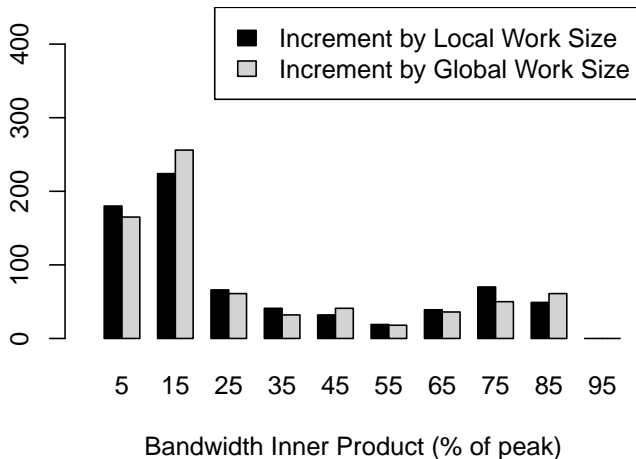
NVIDIA: GTX 285, Tesla K20m



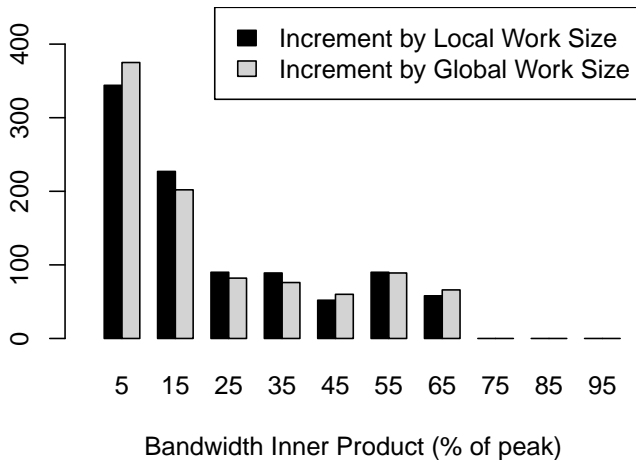
Histograms



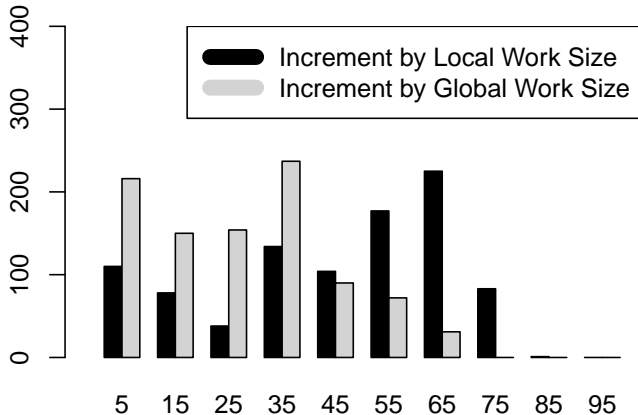
AMD Radeon HD 5850



NVIDIA Tesla K20m



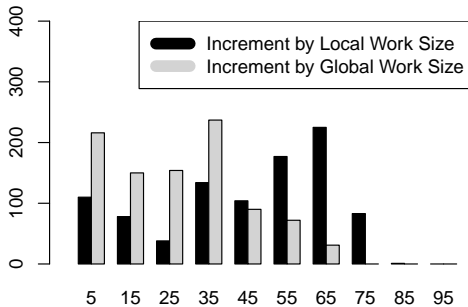
Intel Xeon E5-2670



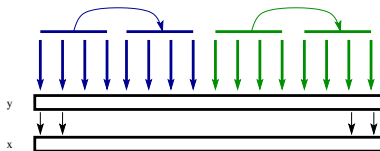
Bandwidth Inner Product (% of theoretical peak)

Benchmark

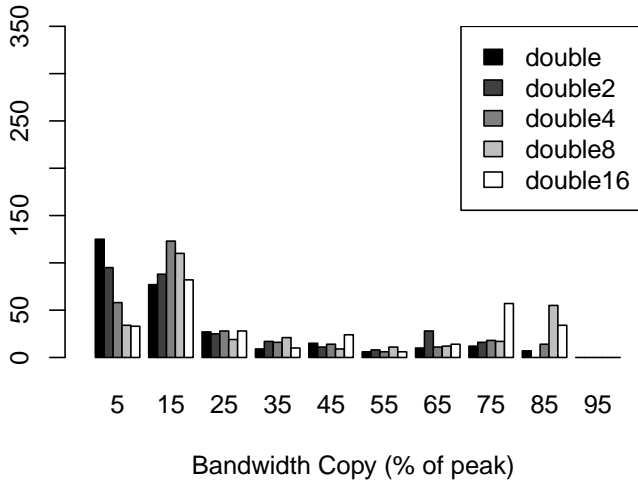
Intel Xeon E5-2670



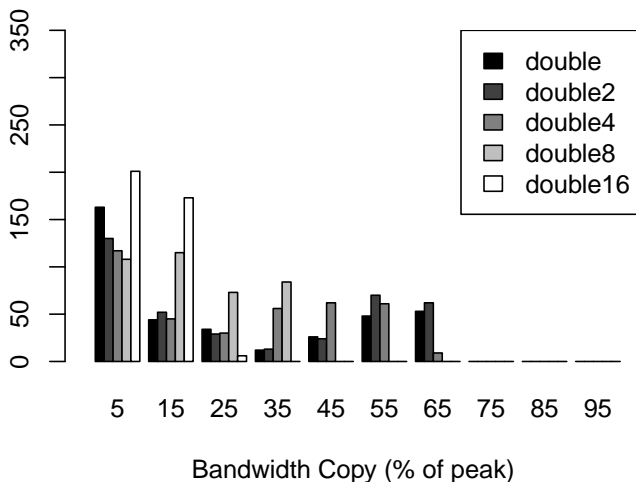
Bandwidth Inner Product (% of theoretical peak)



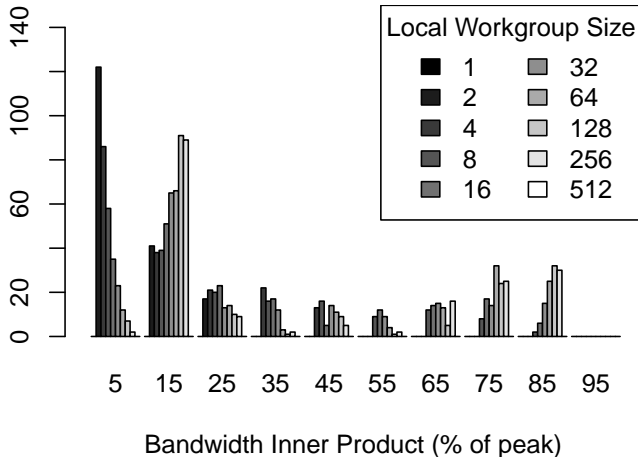
AMD Radeon HD 5850

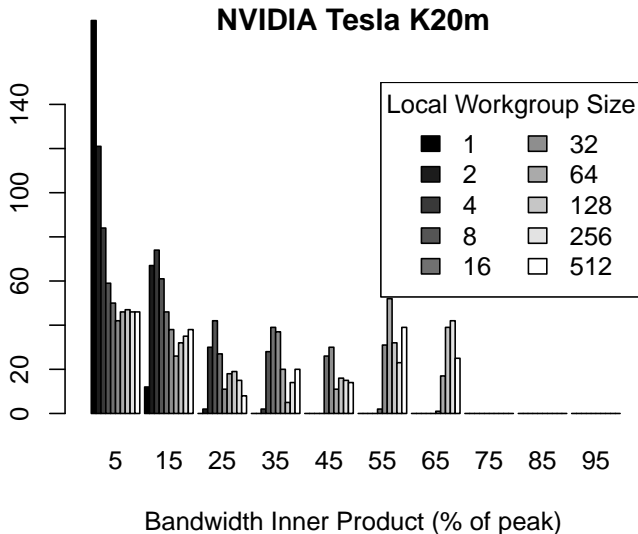


NVIDIA Tesla K20m

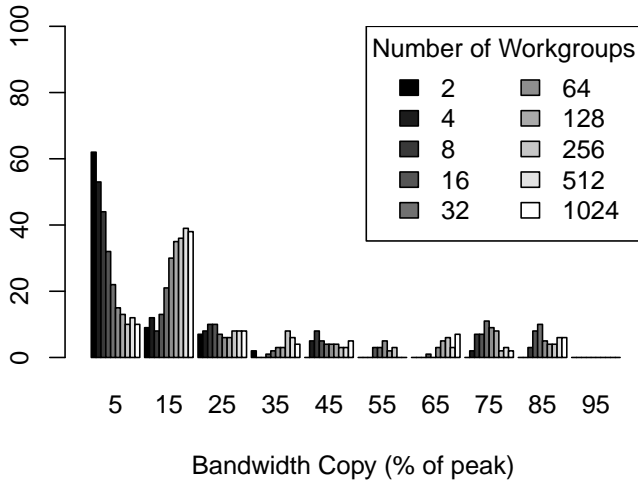


AMD Radeon HD 5850

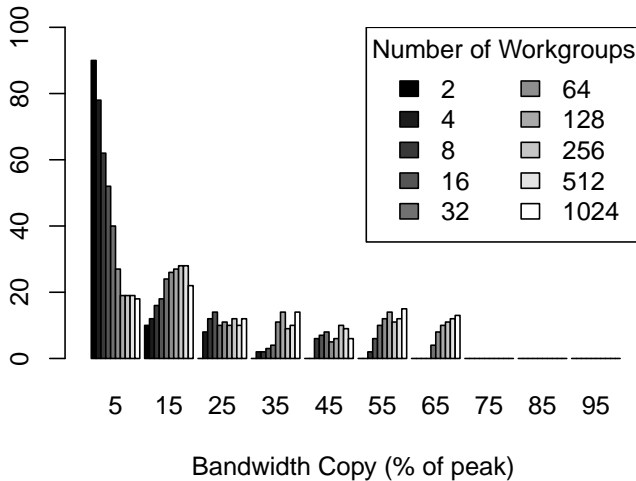




AMD Radeon HD 5850



NVIDIA Tesla K20m

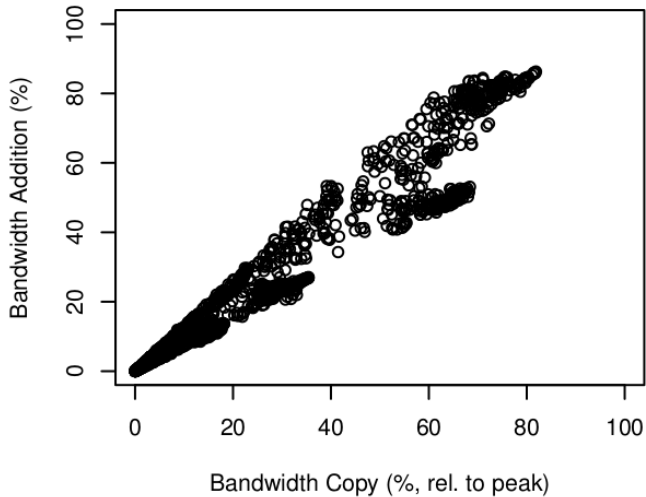


[Addition|Inner Product|Matrix-Vector] vs. Copy Kernel

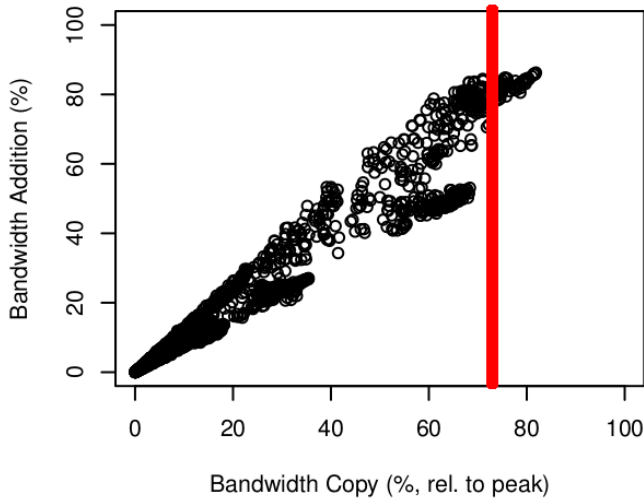
Same Device



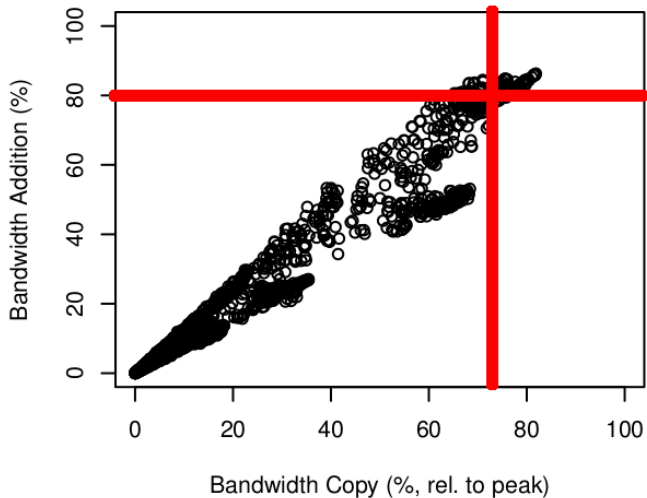
NVIDIA GeForce GTX 285



NVIDIA GeForce GTX 285

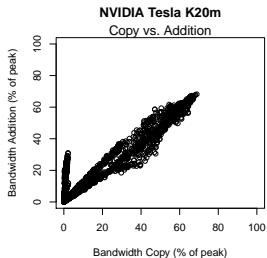


NVIDIA GeForce GTX 285

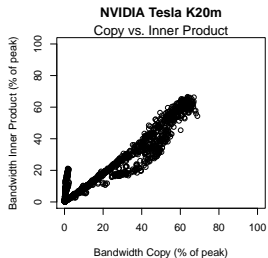


NVIDIA Tesla K20m

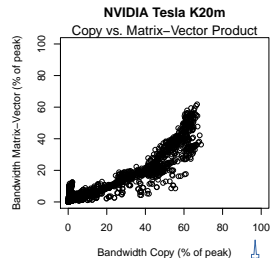
Addition



Inner Product

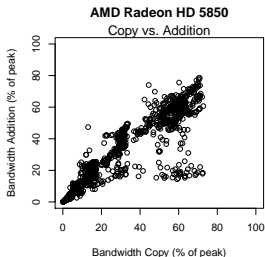


Mat-Vec Product

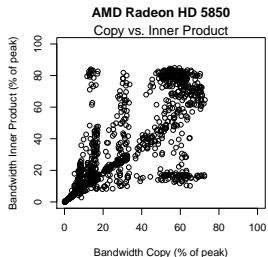


AMD Radeon HD 5850

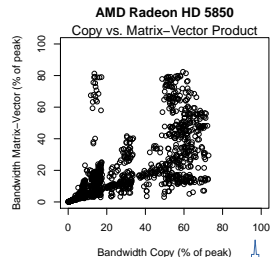
Addition



Inner Product

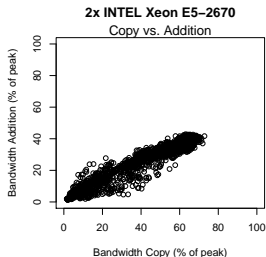


Mat-Vec Product

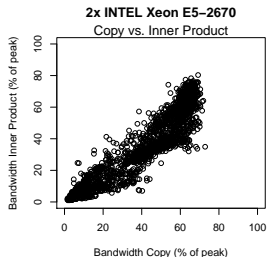


INTEL Dual Xeon E5-2670

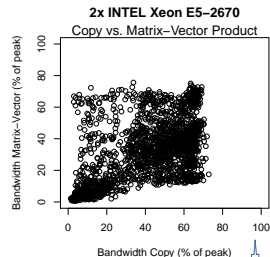
Addition



Inner Product

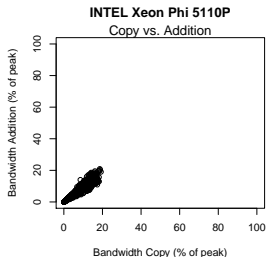


Mat-Vec Product

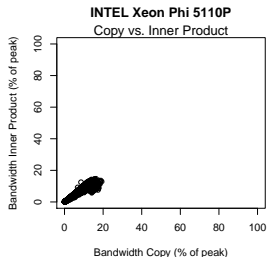


INTEL Xeon Phi

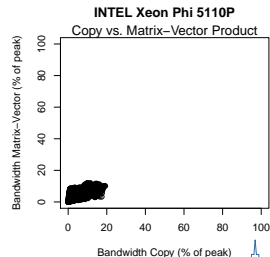
Addition



Inner Product



Mat-Vec Product



Conclusio:

Focus on fastest configurations for copy-kernel sufficient

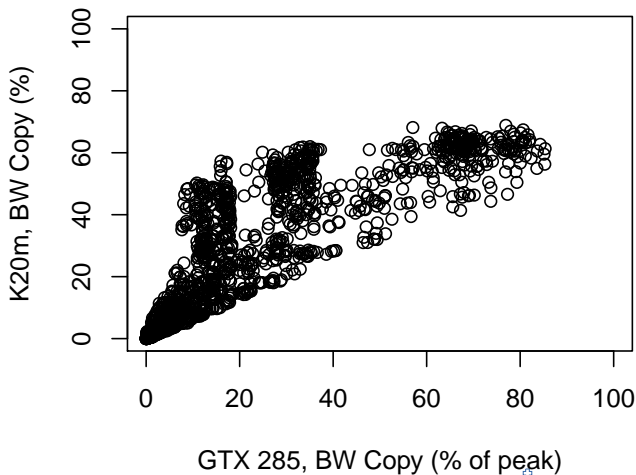


[Copy|Addition|Inner Product|Matrix-Vector] vs. Copy Kernel

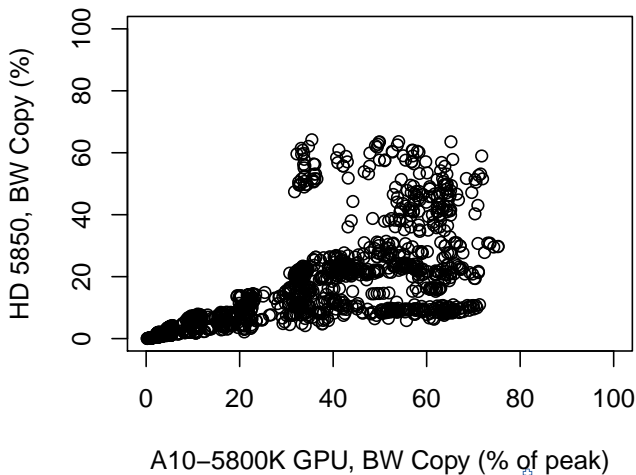
Different Device, Same Vendor



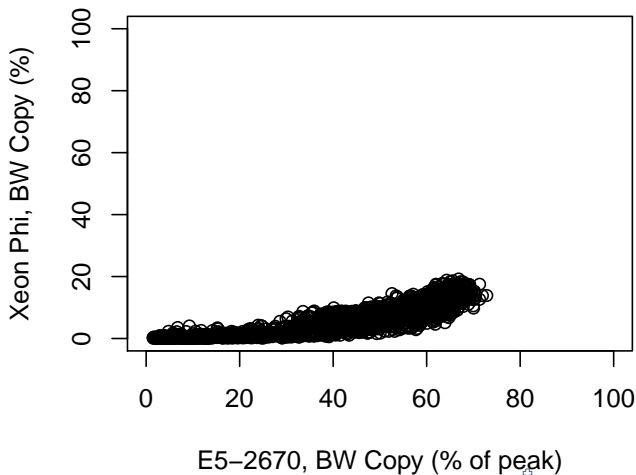
NVIDIA Hardware (x: GTX 285, y: K20m)



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

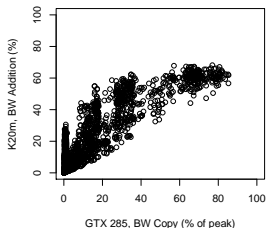


INTEL Hardware (x: Xeon Phi, E5-2670)

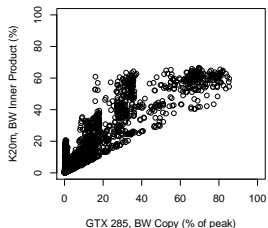


NVIDIA Hardware (x: GTX 285, y: K20m)

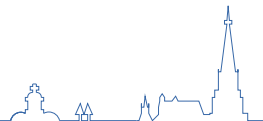
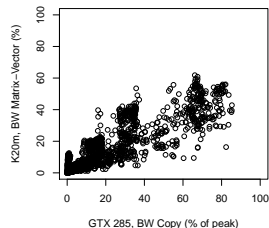
Addition



Inner Product

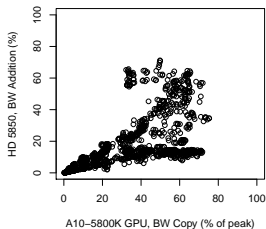


Mat-Vec Product

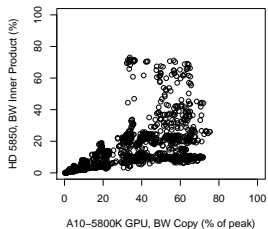


AMD Hardware (x: A10-5800K GPU, y: HD 5850)

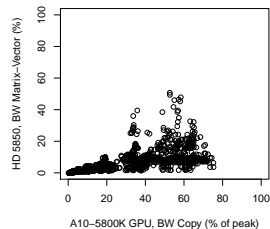
Addition



Inner Product



Mat-Vec Product



Conclusio:

Certain Performance Portability per Vendor



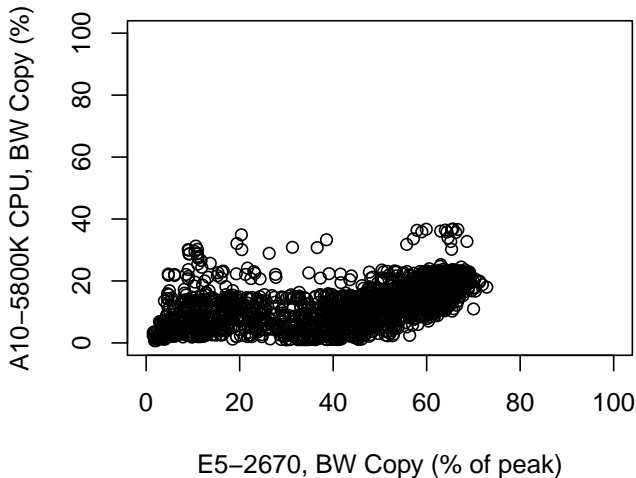
[Copy|Addition|Inner Product|Matrix-Vector] vs. Copy Kernel

Different Device, Different Vendor



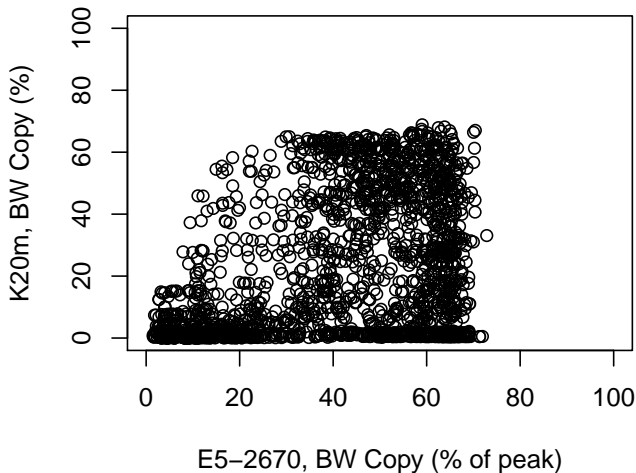
Benchmark

x: INTEL CPU, y: AMD CPU

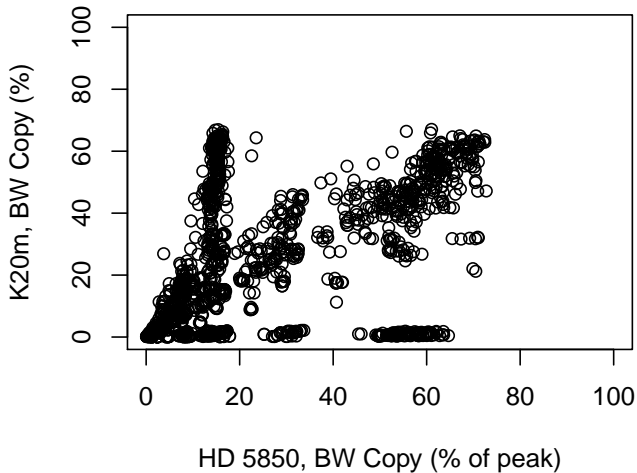


Benchmark

x: INTEL CPU, y: NVIDIA GPU

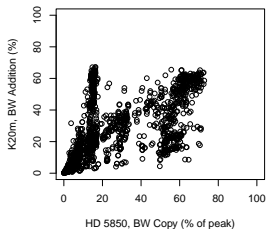


x: AMD GPU, y: NVIDIA GPU

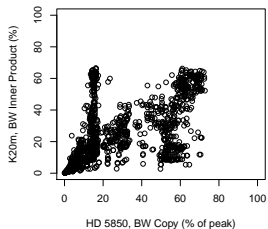


x: AMD HD 5850, y: NVIDIA K20m

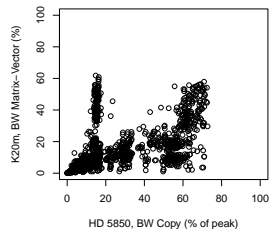
Addition



Inner Product



Mat-Vec Product



Conclusio:

Fast Configurations Across Vendors Exist



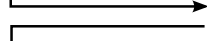
Introduction



Conjugate Gradient - Standard Formulation



Conjugate Gradient - Kernel Parameters



Parameter Study for Portable Performance



Conjugate Gradient - Pipelining



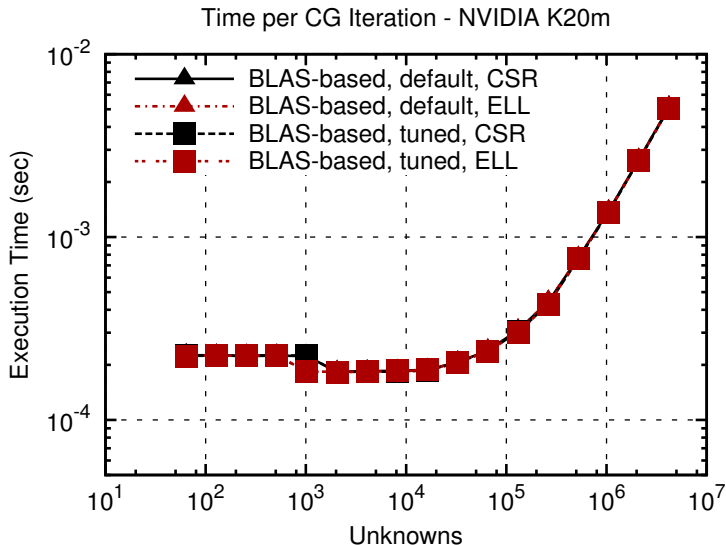
BiCGStab and GMRES - Benchmark Results



Conclusion

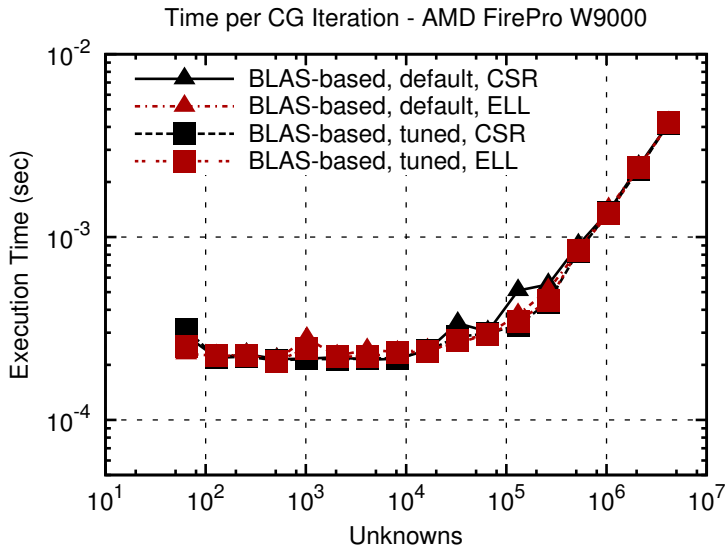


Conjugate Gradients



(2D Finite Difference Discretization)

Conjugate Gradients



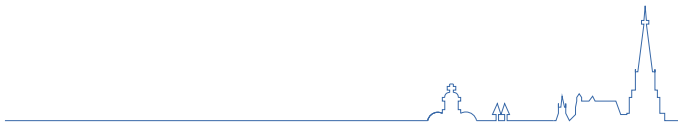
(2D Finite Difference Discretization)

Optimization 3: Rearrange the algorithm

- Remove unnecessary reads

- Remove unnecessary synchronizations

- Use custom kernels instead of standard BLAS



Standard CG

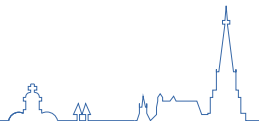
Choose x_0

$$p_0 = r_0 = b - Ax_0$$

For $i = 0$ until convergence

1. Compute and store Ap_i
2. Compute $\langle p_i, Ap_i \rangle$
3. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
4. $x_{i+1} = x_i + \alpha_i p_i$
5. $r_{i+1} = r_i - \alpha_i Ap_i$
6. Compute $\langle r_{i+1}, r_{i+1} \rangle$
7. $\beta_i = \langle r_{i+1}, r_{i+1} \rangle / \langle r_i, r_i \rangle$
8. $p_{i+1} = r_{i+1} + \beta_i p_i$

EndFor



Standard CG

Choose x_0

$$p_0 = r_0 = b - Ax_0$$

For $i = 0$ until convergence

1. Compute and store Ap_i
2. Compute $\langle p_i, Ap_i \rangle$
3. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
4. $x_{i+1} = x_i + \alpha_i p_i$
5. $r_{i+1} = r_i - \alpha_i Ap_i$
6. Compute $\langle r_{i+1}, r_{i+1} \rangle$
7. $\beta_i = \langle r_{i+1}, r_{i+1} \rangle / \langle r_i, r_i \rangle$
8. $p_{i+1} = r_{i+1} + \beta_i p_i$

EndFor

Pipelined CG

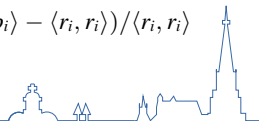
Choose x_0

$$p_0 = r_0 = b - Ax_0$$

For $i = 1$ until convergence

1. $i = 1$: Compute α_0, β_0, Ap_0
2. $x_i = x_{i-1} + \alpha_{i-1} p_{i-1}$
3. $r_i = r_{i-1} - \alpha_{i-1} Ap_i$
4. $p_i = r_i + \beta_{i-1} p_{i-1}$
5. Compute and store Ap_i
6. Compute $\langle Ap_i, Ap_i \rangle, \langle p_i, Ap_i \rangle, \langle r_i, r_i \rangle$
7. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
8. $\beta_i = (\alpha_i^2 \langle Ap_i, Ap_i \rangle - \langle r_i, r_i \rangle) / \langle r_i, r_i \rangle$

EndFor



Standard CG

Choose x_0

$$p_0 = r_0 = b - Ax_0$$

For $i = 0$ until convergence

1. Compute and store Ap_i
2. Compute $\langle p_i, Ap_i \rangle$
3. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
4. $x_{i+1} = x_i + \alpha_i p_i$
5. $r_{i+1} = r_i - \alpha_i Ap_i$
6. Compute $\langle r_{i+1}, r_{i+1} \rangle$
7. $\beta_i = \langle r_{i+1}, r_{i+1} \rangle / \langle r_i, r_i \rangle$
8. $p_{i+1} = r_{i+1} + \beta_i p_i$

EndFor

Pipelined CG

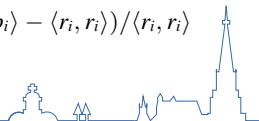
Choose x_0

$$p_0 = r_0 = b - Ax_0$$

For $i = 1$ until convergence

1. $i = 1$: Compute α_0, β_0, Ap_0
2. $x_i = x_{i-1} + \alpha_{i-1} p_{i-1}$
3. $r_i = r_{i-1} - \alpha_{i-1} Ap_i$
4. $p_i = r_i + \beta_{i-1} p_{i-1}$
5. Compute and store Ap_i
6. Compute $\langle Ap_i, Ap_i \rangle, \langle p_i, Ap_i \rangle, \langle r_i, r_i \rangle$
7. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
8. $\beta_i = (\alpha_i^2 \langle Ap_i, Ap_i \rangle - \langle r_i, r_i \rangle) / \langle r_i, r_i \rangle$

EndFor



Standard CG

Choose x_0

$$p_0 = r_0 = b - Ax_0$$

For $i = 0$ until convergence

1. Compute and store Ap_i
2. Compute $\langle p_i, Ap_i \rangle$
3. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
4. $x_{i+1} = x_i + \alpha_i p_i$
5. $r_{i+1} = r_i - \alpha_i Ap_i$
6. Compute $\langle r_{i+1}, r_{i+1} \rangle$
7. $\beta_i = \langle r_{i+1}, r_{i+1} \rangle / \langle r_i, r_i \rangle$
8. $p_{i+1} = r_{i+1} + \beta_i p_i$

EndFor

Pipelined CG

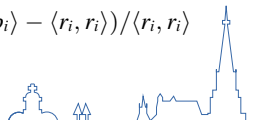
Choose x_0

$$p_0 = r_0 = b - Ax_0$$

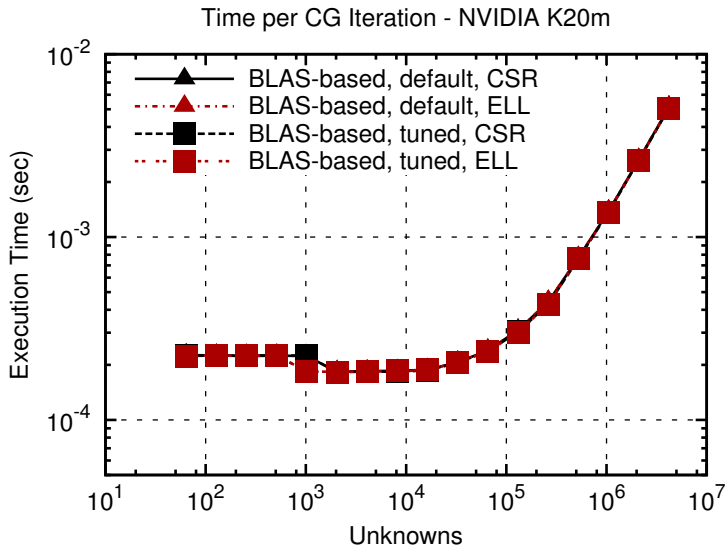
For $i = 1$ until convergence

1. $i = 1$: Compute α_0, β_0, Ap_0
2. $x_i = x_{i-1} + \alpha_{i-1} p_{i-1}$
3. $r_i = r_{i-1} - \alpha_{i-1} Ap_{i-1}$
4. $p_i = r_i + \beta_{i-1} p_{i-1}$
5. Compute and store Ap_i
6. Compute $\langle Ap_i, Ap_i \rangle, \langle p_i, Ap_i \rangle, \langle r_i, r_i \rangle$
7. $\alpha_i = \langle r_i, r_i \rangle / \langle p_i, Ap_i \rangle$
8. $\beta_i = (\alpha_i^2 \langle Ap_i, Ap_i \rangle - \langle r_i, r_i \rangle) / \langle r_i, r_i \rangle$

EndFor

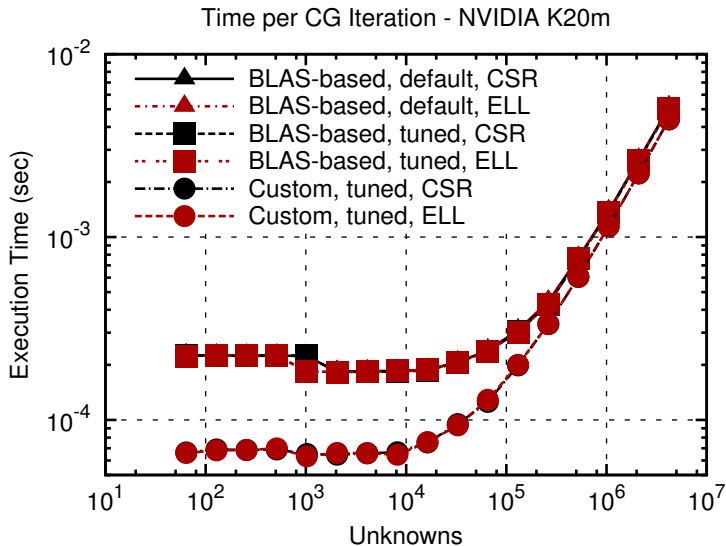


Conjugate Gradients



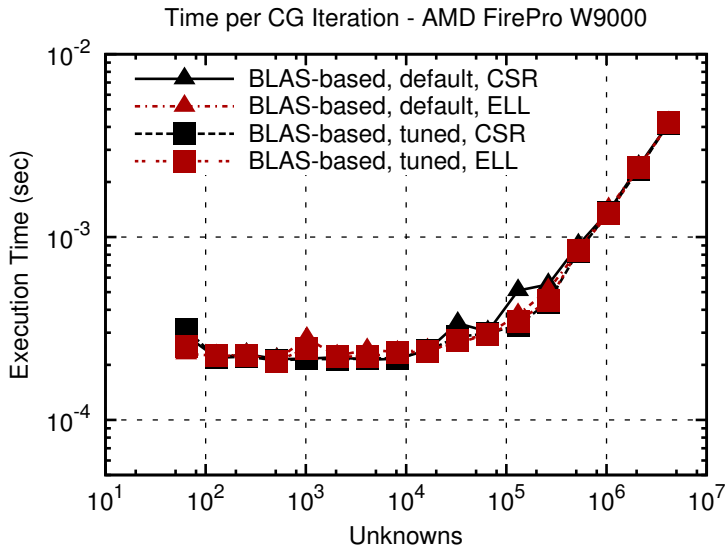
(2D Finite Difference Discretization)

Conjugate Gradients



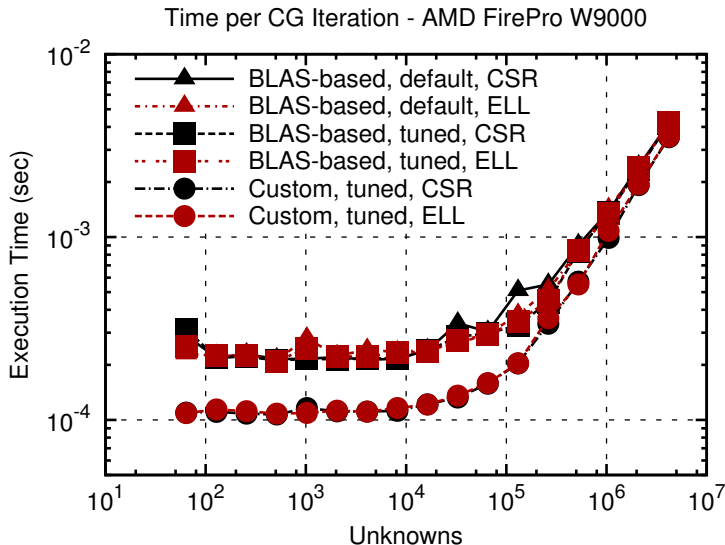
(2D Finite Difference Discretization)

Conjugate Gradients



(2D Finite Difference Discretization)

Conjugate Gradients



(2D Finite Difference Discretization)

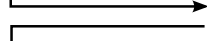
Introduction



Conjugate Gradient - Standard Formulation



Conjugate Gradient - Kernel Parameters



Parameter Study for Portable Performance



Conjugate Gradient - Pipelining



BiCGStab and GMRES - Benchmark Results



Conclusion



BiCGStab

Similar to CG

Two SpMV per iteration

Pipelining: 4 kernel launches instead of 12

GMRES

Store Krylov basis

Orthonormalization in each step

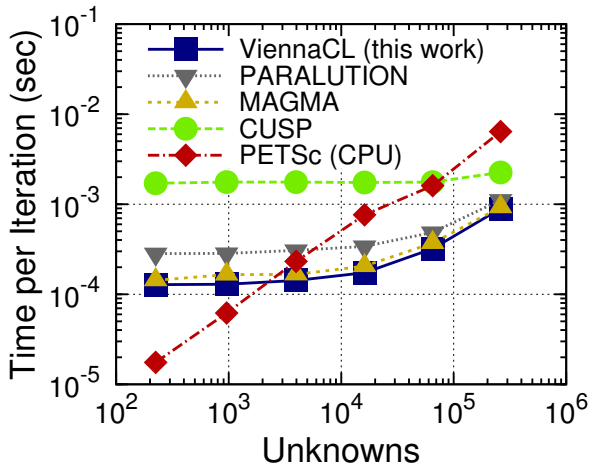
Pipelining: 3 kernel launches

Benchmark Setup

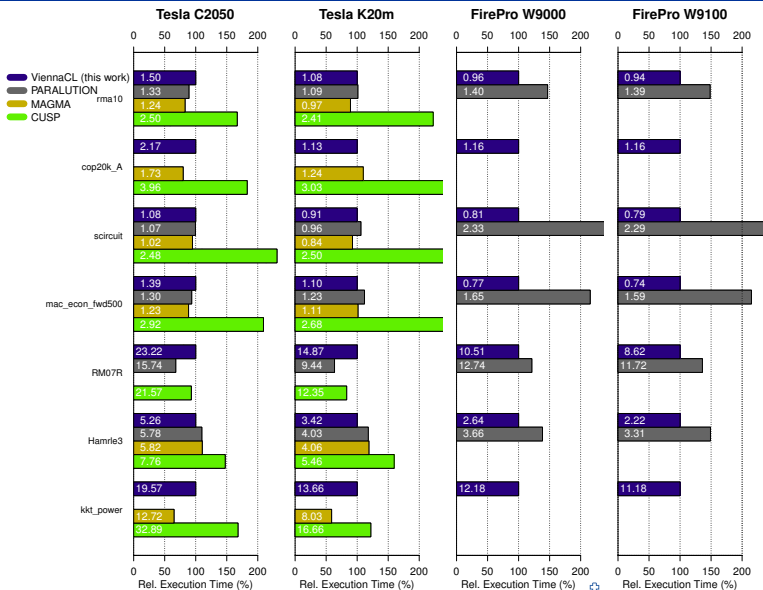
Poisson equation in 2D

GPUs from NVIDIA and AMD

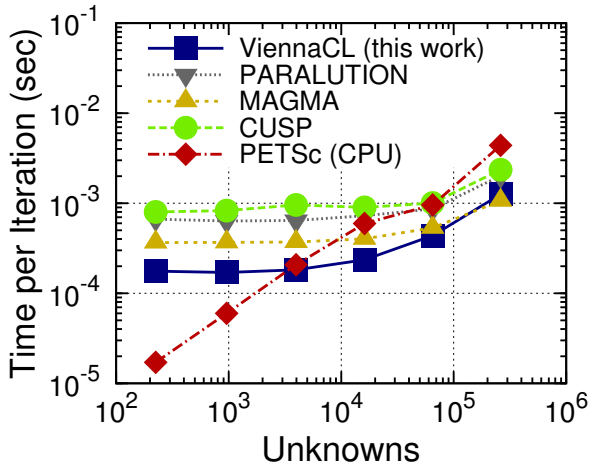




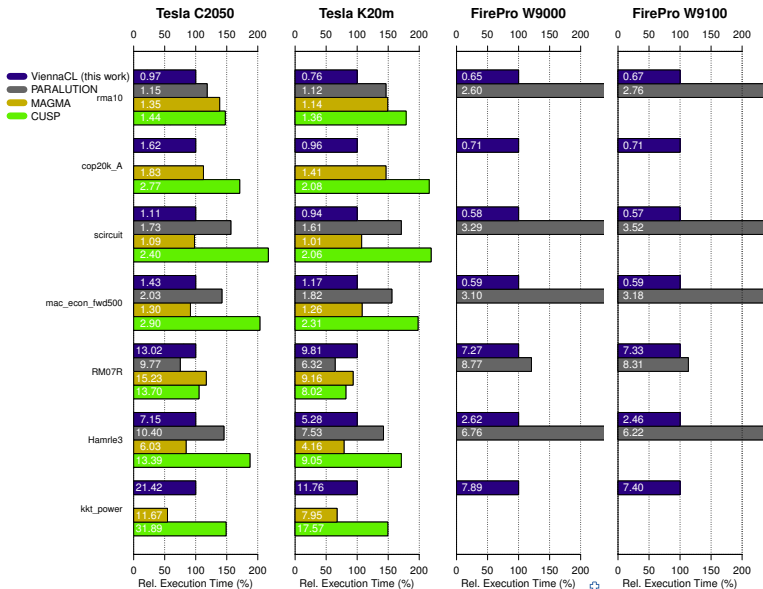
BiCGStab Benchmarks



GMRES Benchmarks



GMRES Benchmarks



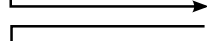
Introduction



Conjugate Gradient - Standard Formulation



Conjugate Gradient - Kernel Parameters



Parameter Study for Portable Performance



Conjugate Gradient - Pipelining



BiCGStab and GMRES - Benchmark Results



Conclusion



BiCGStab

Similar to CG

Two SpMV per iteration

Pipelining: 4 kernel launches instead of 12

GMRES

Store Krylov basis

Orthonormalization in each step

Pipelining: 3 kernel launches

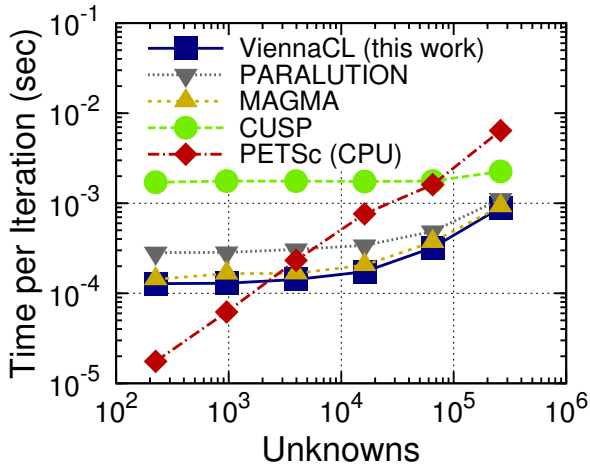
Benchmark Setup

Poisson equation in 2D

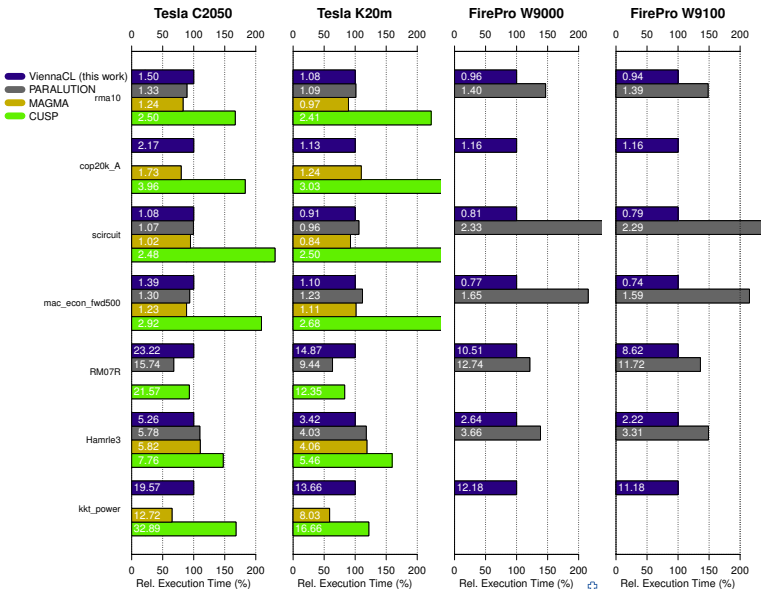
GPUs from NVIDIA and AMD



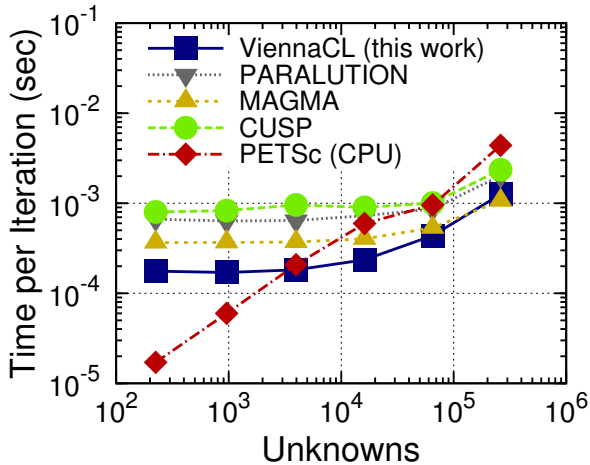
BiCGStab Benchmarks



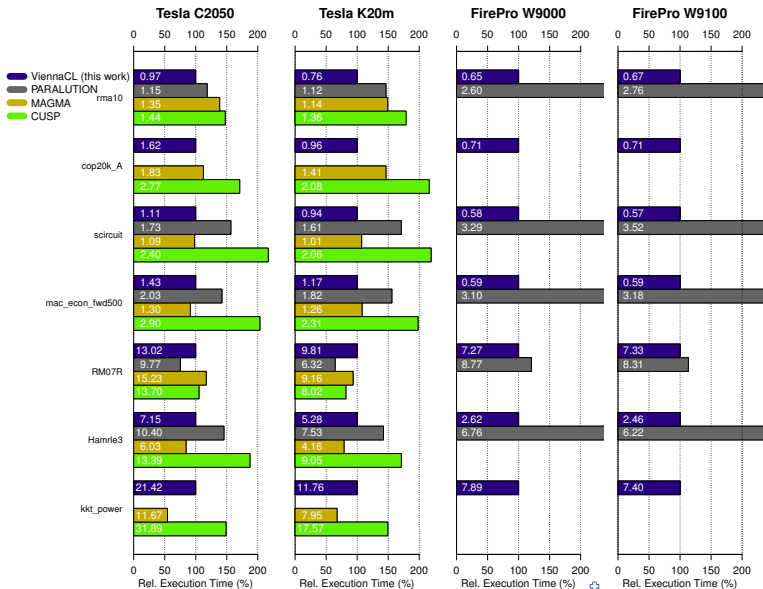
BiCGStab Benchmarks



GMRES Benchmarks



GMRES Benchmarks



Performance-Portable Code for GPUs

- Start with 128 work items

- Refrain from using vector datatypes

- Let each workgroup work on a contiguous piece of memory

Pipelined Iterative Solvers

- Reduced number of kernel launches

- On-chip data reuse

- Faster than BLAS-based implementations

Best Practices for GPU Computing

- FLOPs are (almost always) for free

- Work on large enough data

- Avoid unnecessary PCI-Express communication

