

# FEM Integration with Quadrature on GPUs



Matthew Knepley<sup>1</sup>, Karl Rupp<sup>2</sup>, Andy R. Terrel<sup>3</sup>

`rupp@iue.tuwien.ac.at`  
`http://karlrupp.net/`



<sup>1</sup> University of Chicago, USA

<sup>2</sup> Vienna University of Technology, Austria

<sup>3</sup> Continuum Analytics, USA

March 18th, 2015

## Finite Element Method

Several basis functions per element

Evaluation of integrals on each element

## General Weak Form

Residual formulation for test function  $\phi$

$$\int_{\Omega} \phi \cdot f_0(u, \nabla u) + \nabla \phi : \mathbf{f}_1(u, \nabla u) = 0.$$

## Examples

Laplace:  $f_0 \equiv 0$ ,  $\mathbf{f}_1 \equiv \nabla u$

Poisson:  $f_0 \equiv g$ ,  $\mathbf{f}_1 \equiv \nabla u$



$$\int_{\Omega} \phi \cdot f_0(u, \nabla u) + \nabla \phi : \mathbf{f}_1(u, \nabla u) = 0.$$

## Element-Wise General Weak Form

Evaluation using quadrature

$$\sum_e \mathcal{E}_e^T \left[ B^T W f_0(u^q, \nabla u^q) + \sum_k D_k^T W \mathbf{f}_1^k(u^q, \nabla u^q) \right] = 0$$

$\mathcal{E}$  ... global vector

$W$  ... quadrature weights

$B, D_k$  ... reduction operations for global basis coefficients

## Parallelization Options

Across elements

Quadrature points

Basis functions



## Parallelization Across Elements

- Large memory per thread

- Synchronizations with neighbor elements

- [Cecka et al. 2011; Taylor et al. 2008; Williams 2012]

## Parallelization per Quadrature Point

- No memory overhead

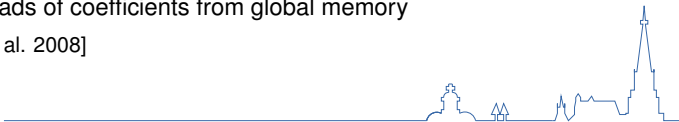
- Too many synchronizations

## Parallelization via Basis Functions

- Very little local memory

- Repeated loads of coefficients from global memory

- [Dabrowski et al. 2008]

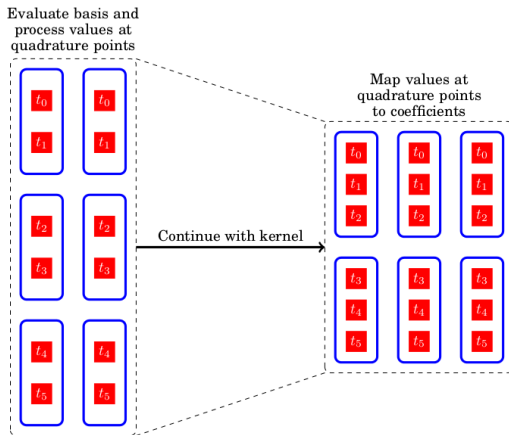


## Thread Block Works on Multiple Elements

Number of quadrature points  $N_q$

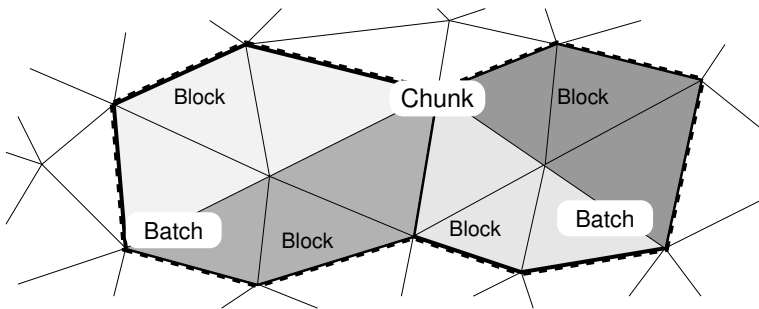
Number of basis functions  $N_b$

Minimum number of elements  $\text{LCM}(N_q, N_b)$



## High Level Decomposition

- Chunks - Cells processed by each thread workgroup
- Batches - Cells processed with one thread transposition
- Blocks - Smallest unit of execution



## OpenCL-enabled Hardware

NVIDIA GTX 470

NVIDIA GTX 580

NVIDIA Tesla K20m

AMD FirePro W9100

(AMD A10-5800K)

## Comparisons

Single vs. double precision

2D vs. 3D

## Invariants

Variable coefficients

First-order FEM

Poisson equation



## Choice of Block and Batch Numbers

NVIDIA GTX 470

Performance in GFLOPs/sec

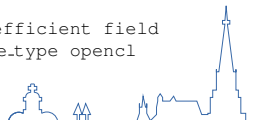
Actual choice not very sensitive

Blocks	Batches					
	16	20	24	28	32	36
4	113	120	118	122	<b>137</b>	119
8	109	116	113	120	108	117
12	102	112	110	109	115	113
16	108	100	99	111	130	106

(2D triangular mesh, variable coefficients, single precision, NVIDIA GTX 470)

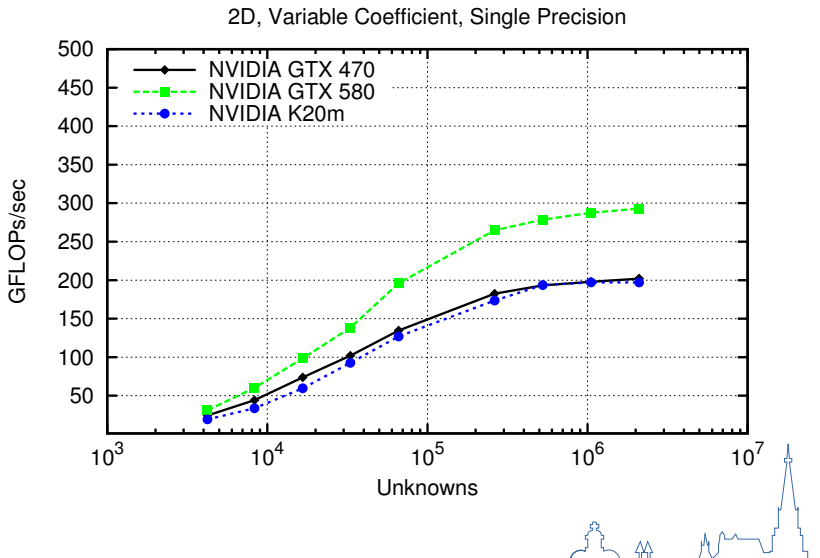
PETSc SNES ex12:

```
./ex12 -petscspace_order 1 -run_type perf -variable_coefficient field  
-refinement_limit 0.00001 -show_solution false -petscfe_type openc1  
-petscfe_num_blocks 4 -petscfe_num_batches 16
```

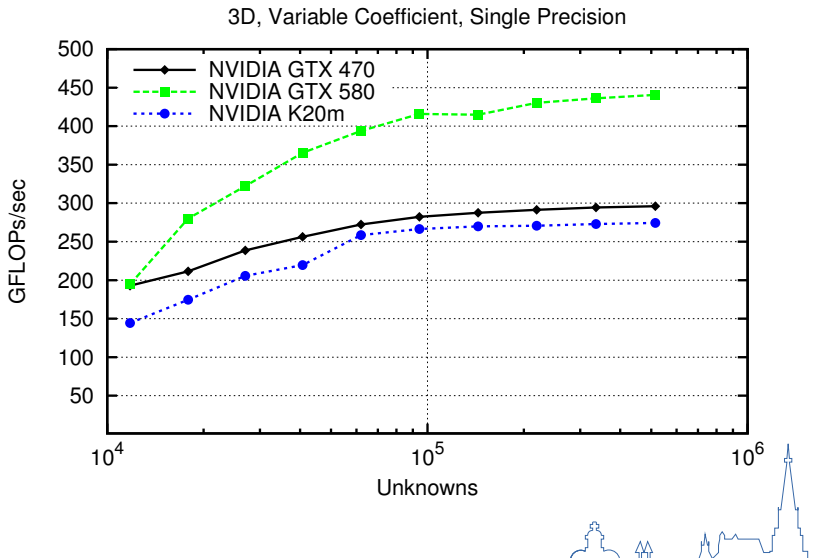




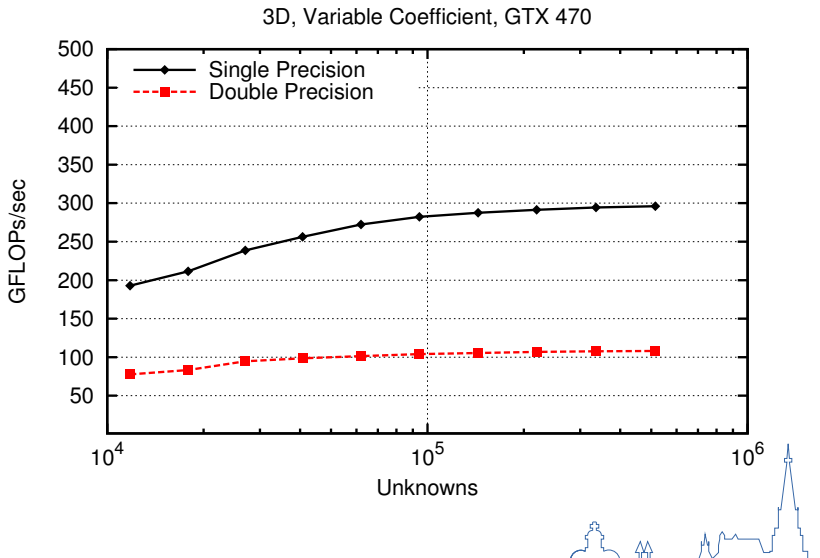
# Benchmark



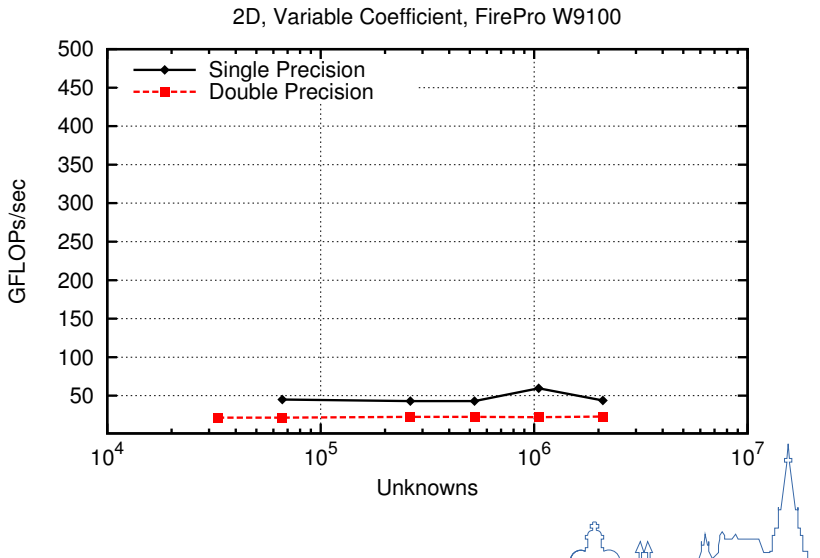
# Benchmark



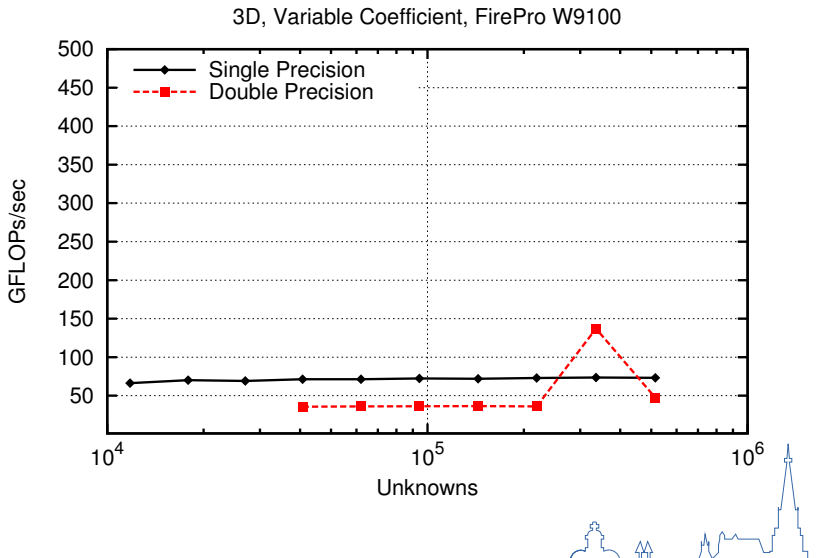
# Benchmark



# Benchmark



# Benchmark



## Limiting Factor?

GTX 470: 134 GB/sec memory bandwidth (theoretical)

GTX 470: 1088 GFLOPs/sec peak (theoretical)

## Arithmetic Intensity

Count FLOPs and bytes loaded/stored

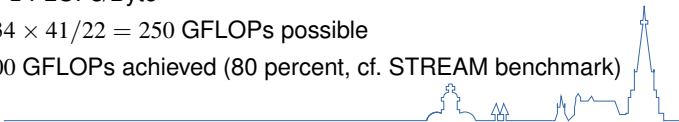
$$\beta = \frac{[(2 + (2 + 2d)d)N_{bt}N_q + 2dN_{comp}N_q + (2 + 2d)dN_qN_{bt}]N_{bs}N_{bl}}{4N_t((d^2 + 1) + N_{bt} + (d + 1)N_q)}$$

## 2D Mesh, First-Order FEM, Single Precision

$$\beta = 41/22 \approx 2 \text{ FLOPs/Byte}$$

GTX 470:  $134 \times 41/22 = 250$  GFLOPs possible

GTX 470: 200 GFLOPs achieved (80 percent, cf. STREAM benchmark)



## FEM Quadrature on GPUs

“Matrix-Free”

Higher arithmetic intensity

## Performance Results

Good performance on NVIDIA GPUs and AMD APUs

5x improvements for discrete AMD GPUs desired

## Performance Modeling

Performance limited by memory bandwidth

Excellent prediction accuracy

## Reproducibility

PETSc, SNES tutorial, ex12

