

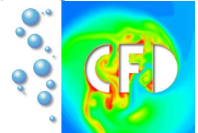
GPU Acceleration of Unmodified CSM and CFD Solvers

Dominik Göddeke
Sven H.M. Buijssen, Hilmar Wobker and Stefan Turek

Angewandte Mathematik und Numerik
TU Dortmund, Germany

`dominik.goeddeke@math.tu-dortmund.de`

High Performance Computing & Simulation
Workshop on Architecture-aware Simulation and Computing
Leipzig, June 22, 2009



Scientific computing is in the middle of a paradigm shift

ILP wall memory wall characteristic feature size
heat power consumption leaking voltage

Hardware evolves towards parallelism and heterogeneity

multicore CPUs Cell BE processor GPUs

Emerging manycore architectures

accelerators algorithm design for 10000s of threads

FEAST – Hardware-oriented Numerics

Fully adaptive grids

Maximum flexibility

'Stochastic' numbering

Unstructured sparse matrices

Indirect addressing, very slow.

Locally structured grids

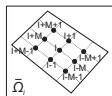
Logical tensor product

Fixed banded matrix structure

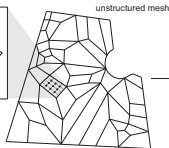
Direct addressing (\Rightarrow fast)

r -adaptivity

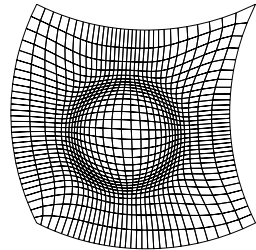
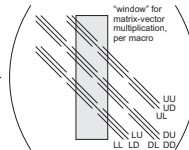
Unstructured macro mesh of tensor product subdomains



hierarchically
refined subdomain
(= "macro"),
rowwise numbered



unstructured mesh



ScaRC – Scalable Recursive Clustering

- Minimal overlap by extended Dirichlet BCs
- Hybrid multilevel domain decomposition method
- Inspired by parallel MG ("best of both worlds")
 - Multiplicative vertically (between levels), global coarse grid problem (MG-like)
 - Additive horizontally: block-Jacobi / Schwarz smoother (DD-like)
- Hide local irregularities by MGs within the Schwarz smoother
- Embed in Krylov to alleviate Block-Jacobi character

global BiCGStab

preconditioned by

global multilevel (V 1+1)

additively smoothed by

for all Ω_i : **local multigrid**

coarse grid solver: UMFPACK

Block-structured systems

- Guiding idea: Tune scalar case once per architecture instead of over and over again per application
- Equation-wise ordering of the unknowns
- Block-wise treatment enables multivariate ScaRC solvers

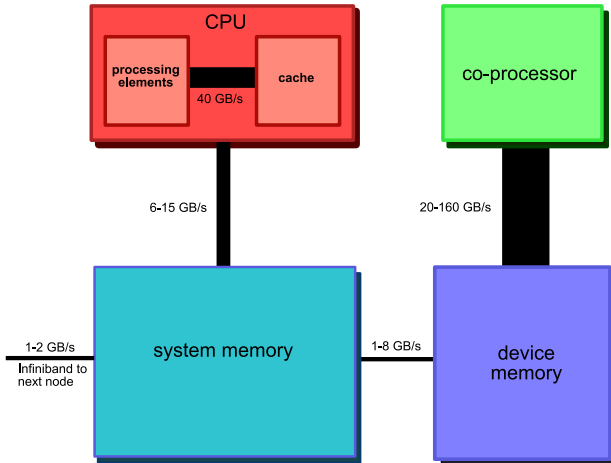
Examples

- Linearised elasticity with compressible material (2x2 blocks)
- Saddle point problems: Stokes, linearised elasticity with (nearly) incompressible material, Navier-Stokes with stabilisation (3x3 blocks, three zero Blocks for Stokes)

$$\begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} = \mathbf{f}, \quad \begin{pmatrix} \mathbf{A}_{11} & \mathbf{0} & \mathbf{B}_1 \\ \mathbf{0} & \mathbf{A}_{22} & \mathbf{B}_2 \\ \mathbf{B}_1^T & \mathbf{B}_2^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{p} \end{pmatrix} = \mathbf{f},$$

\mathbf{A}_{11} and \mathbf{A}_{22} correspond to scalar (elliptic) operators
 \Rightarrow Tuned linear algebra **and** tuned solvers

Co-processor integration into FEAST



Level	Core2Duo (double)		GTX 280 (mixed)		
	time(s)	MFLOP/s	time(s)	MFLOP/s	speedup
7	0.021	1405	0.009	2788	2.3x
8	0.094	1114	0.012	8086	7.8x
9	0.453	886	0.026	15179	17.4x
10	1.962	805	0.073	21406	26.9x

- Poisson on unitsquare, Dirichlet BCs, TP grid, *not a matrix stencil*
- 1M DOF, multigrid, FE-accurate in less than 0.1 seconds!
- Converges against wrong solution in single precision
- 27x faster than CPU, exactly same results as pure double
- 1.7x faster than pure double on GPU
- 8800 GTX (correction loop on CPU): 0.44 seconds on level 10

global BiCGStab

preconditioned by

global multilevel (V 1+1)

additively smoothed by

for all Ω_i : **local multigrid**

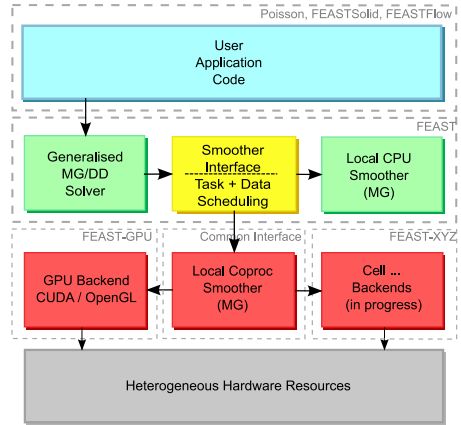
coarse grid solver: UMFPACK

All outer work: CPU, double

Local MGs: GPU, single

GPU performs preconditioning

Applicable to many co-
processors



General approach

- Balance acceleration potential and integration effort
- Accelerate many different applications built on top of one central FE and solver toolkit
- Diverge code paths as late as possible
- No changes to application code!
- Retain all functionality
- Do not sacrifice accuracy

Challenges

- Heterogeneous task assignment to maximise throughput
- Limited device memory (modeled as huge L3 cache)
- Overlapping CPU and GPU computations
- Building dense accelerated clusters

Some results

$$\begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} = \mathbf{f}$$

$$\begin{pmatrix} (2\mu + \lambda)\partial_{xx} + \mu\partial_{yy} & (\mu + \lambda)\partial_{xy} \\ (\mu + \lambda)\partial_{yx} & \mu\partial_{xx} + (2\mu + \lambda)\partial_{yy} \end{pmatrix}$$

global multivariate BiCGStab

block-preconditioned by

Global multivariate multilevel (V 1+1)

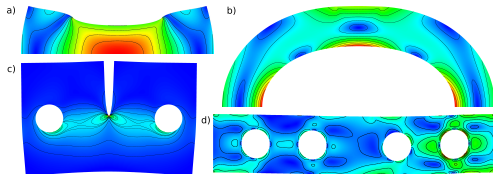
additively smoothed (block GS) by

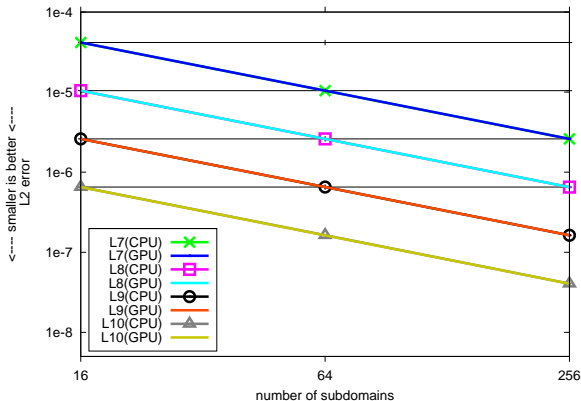
for all Ω_i : solve $\mathbf{A}_{11}\mathbf{c}_1 = \mathbf{d}_1$ by
local scalar multigrid

update RHS: $\mathbf{d}_2 = \mathbf{d}_2 - \mathbf{A}_{21}\mathbf{c}_1$

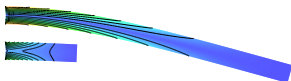
for all Ω_i : solve $\mathbf{A}_{22}\mathbf{c}_2 = \mathbf{d}_2$ by
local scalar multigrid

coarse grid solver: UMFPACK

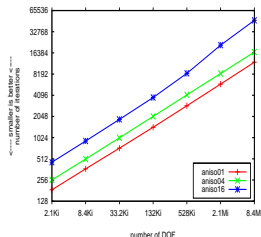




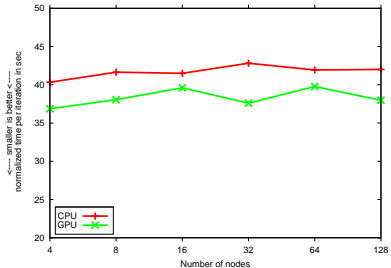
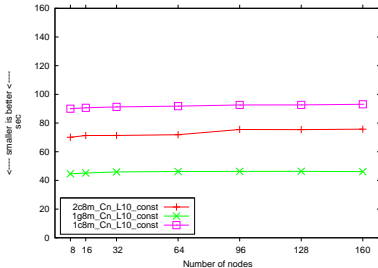
- Same results for CPU and GPU
- L_2 error against analytically prescribed displacements
- Tests on 32 nodes, 512 M DOF



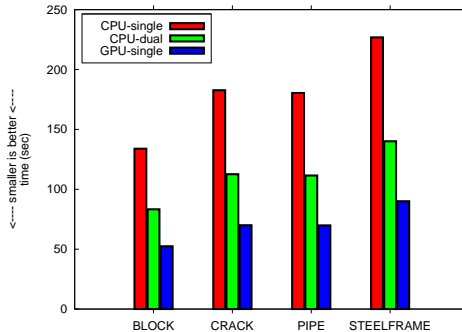
Cantilever beam, aniso 1:1, 1:4, 1:16
Hard, very ill-conditioned CSM test
CG solver: $> 2x$ iterations per refinement
GPU-ScaRC solver: same results as CPU



aniso04 refinement L	Iterations		Volume		y-Displacement	
	CPU	GPU	CPU	GPU	CPU	GPU
8	4	4	1.6087641E-3	1.6087641E-3	-2.8083499E-3	-2.8083499E-3
9	4	4	1.6087641E-3	1.6087641E-3	-2.8083628E-3	-2.8083628E-3
10	4.5	4.5	1.6087641E-3	1.6087641E-3	-2.8083667E-3	-2.8083667E-3
aniso16						
8	6	6	6.7176398E-3	6.7176398E-3	-6.6216232E-2	-6.6216232E-2
9	6	5.5	6.7176427E-3	6.7176427E-3	-6.6216551E-2	-6.6216552E-2
10	5.5	5.5	6.7176516E-3	6.7176516E-3	-6.6217501E-2	-6.6217502E-2



- Outdated cluster, dual Xeon EM64T,
- one NVIDIA Quadro FX 1400 per node (one generation behind the Xeons, 20 GB/s BW)
- Poisson problem (left): up to 1.3 B DOF, 160 nodes
- Elasticity (right): up to 1 B DOF, 128 nodes

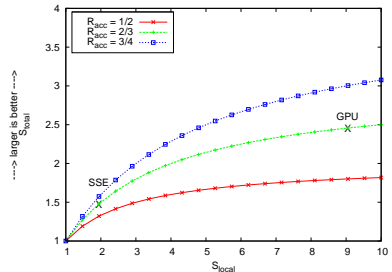


- 16 nodes, Opteron X2 2214,
- NVIDIA Quadro FX 5600 (76 GB/s BW), OpenGL
- Problem size 128 M DOF
- Dualcore 1.6x faster than singlecore
- GPU 2.6x faster than singlecore, 1.6x than dual

Speedup analysis

- Addition of GPUs increases resources
- \Rightarrow Correct model: strong scalability inside each node
- Accelerable fraction of the elasticity solver: $2/3$
- Remaining time spent in MPI and the outer solver

Accelerable fraction R_{acc} : 66%
Local speedup S_{local} : 9x
Total speedup S_{total} : 2.6x
Theoretical limit S_{max} : 3x



$$\begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} & \mathbf{B}_1 \\ \mathbf{A}_{21} & \mathbf{A}_{22} & \mathbf{B}_2 \\ \mathbf{B}_1^T & \mathbf{B}_2^T & \mathbf{C} \end{pmatrix} \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \mathbf{p} \end{pmatrix} = \begin{pmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \\ \mathbf{g} \end{pmatrix}$$

- 4-node cluster
- Opteron X2 2214
- GeForce 8800 GTX (90 GB/s BW), CUDA
- Driven cavity and channel flow around a cylinder

fixed point iteration

solving linearised subproblems with

global BiCGStab (reduce initial residual by 1 digit)
Block-Schurcomplement preconditioner

1) approx. solve for velocities with

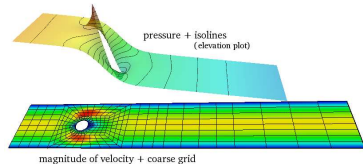
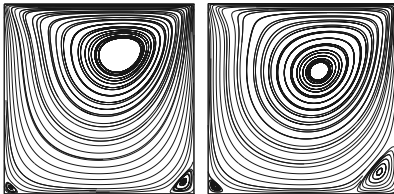
global MG (V 1+0), additively smoothed by

for all Ω_i : solve for \mathbf{u}_1 with
local MG

for all Ω_i : solve for \mathbf{u}_2 with
local MG

2) update RHS: $\mathbf{d}_3 = -\mathbf{d}_3 + \mathbf{B}^T(\mathbf{c}_1, \mathbf{c}_2)^T$

3) scale $\mathbf{c}_3 = (\mathbf{M}_p^L)^{-1} \mathbf{d}_3$



Speedup analysis

	R_{acc}		S_{local}		S_{total}	
	L9	L10	L9	L10	L9	L10
DC Re100	41%	46%	6x	12x	1.4x	1.8x
DC Re250	56%	58%	5.5x	11.5x	1.9x	2.1x
Channel flow	60%	–	6x	–	1.9x	–

Important consequence: Ratio between assembly and linear solve changes significantly

DC Re100		DC Re250		Channel flow	
plain	accel.	plain	accel.	plain	accel.
29:71	50:48	11:89	25:75	13:87	26:74

Conclusions

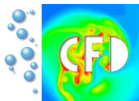
- Hardware-oriented numerics prevents existing codes being worthless in a few years
- Mixed precision schemes exploit the available bandwidth without sacrificing accuracy
- GPUs as local preconditioners in a large-scale parallel FEM package
- Not limited to GPUs, applicable to all kinds of hardware accelerators
- Minimally invasive approach, no changes to application code
- Excellent local acceleration, global acceleration limited by 'sequential' part
- Future work: Design solver schemes with higher acceleration potential without sacrificing numerical efficiency

Collaborative work with

FEAST group (TU Dortmund)

Robert Strzodka (Max Planck Institut Informatik)

Jamaludin Mohd-Yusof, Patrick McCormick (Los Alamos National Laboratory)



<http://www.mathematik.tu-dortmund.de/~goeddeke>

Supported by Deutsche Forschungsgemeinschaft, project TU 102/22-1, TU 102/22-2, TU 102/27-1, TU102/11-3