

Collaboration

- Dominik Göddeke, Stefan Turek, FEAST Group (Dortmund University of Technology)
- Jamaludin Mohd-Yusof, Patrick McCormick (Advanced Computing Lab, LANL)
- Robert Strzodka, Max Planck Center (Max Planck Institut Informatik)







What is a Mixed Precision Method?

- Definition: A method that uses different precisions in its computations
- Example: double(a) + double(float(b) + float(c))
- Typical usage: Mix of single and double precision floating point computations

 Goal: Obtain the same accuracy but better performance with more low precision computations

Mixed Precision Performance Gains

Bandwidth bound algorithm

- 64 bit = 1 double = 2 floats
- More variables per bandwidth (comp. intensity up)
- More variables per storage (data block size up)
- Applies to all memory levels:
 disc → main → device → local → register

Computation bound algorithm

- 1 double multiplier ≈ 4 float multiplier (quadratic)
- -1 double adder ≈ 2 float adder (linear)
- Multipliers are much bigger than adders
 - → Quadrupled computational efficiency

Overview

Why Bother with Mixed Precision?

Precision and Accuracy

Floating Point Operations

Mixed Precision Iterative Refinement



Roundoff and Cancellation

Roundoff examples for the float s23e8 format

```
additive roundoff a=1+0.00000004=1.000000004=_{fl}1 multiplicative roundoff b=1.0002*0.9998=0.999999996=_{fl}1 cancellation c \in \{a,b\} (c-1)*10^8=\pm 4=_{fl}0
```

Cancellation promotes the small error 0.00000004 to the absolute error 4 and a relative error of order one.

Order of operations can be crucial:

```
1 + 0.00000004 - 1 =_{fl} 0

1 - 1 + 0.00000004 =_{fl} 0.00000004
```

With the double s52e11 format no problems above, but ...

An Instructive Example

Evaluating f(x,y) with powers as multiplications [S.M. Rump, 1988]

$$f(x,y) = (333.75 - x^2) y^6 + x^2 (11x^2y^2 - 121y^4 - 2) + 5.5y^8 + x/(2y)$$

for
$$x_0 = 77617$$
, $y_0 = 33096$ gives

float s23e8 1.1726

double s52e11 1.17260394005318

long double s63e15 1.172603940053178631

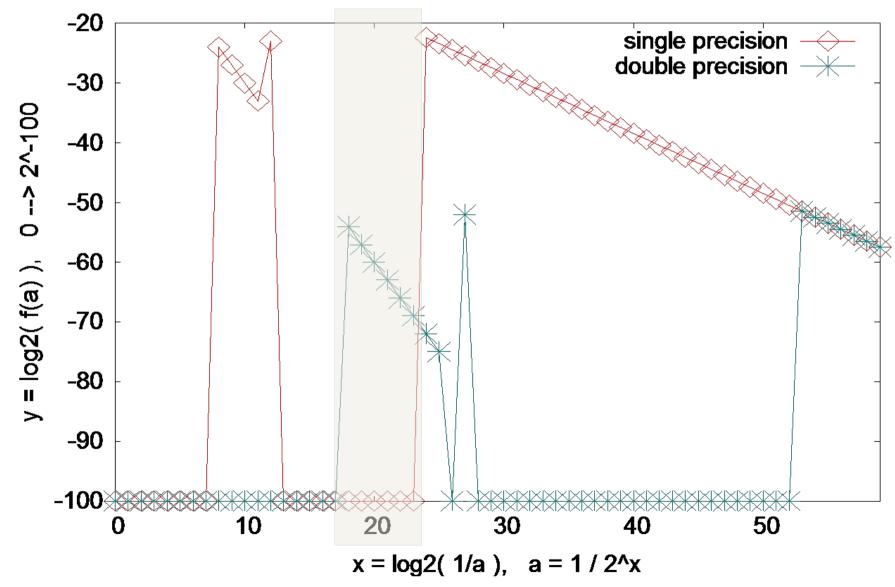
This is all wrong, even the sign is wrong!! The correct result is

-0.82739605994682136814116509547981629...

Lesson learnt: Computational Precision ≠ Accuracy of Result

The Erratic Roundoff Error

Roundoff error for: $0 = f(a) = |(1+a)^3 - (1+3a^2) - (3a+a^3)|$



The Dominant Data Error

- Data error occurs when the exact value has to be truncated for storage in the binary format, e.g.
 - $-\pi$, $\sqrt{2}$, $\sin(2)$, $\exp(2)$, 1/3, ...
 - In fact, any value, e.g. 0.1, except combinations of 2^b

- So more precision is usually better because
 - for float s23e8: $1 + 4e-8 =_{fl} 1$
 - for double s52e11: $1 + 4e-15 =_{fl} 1$
- How can float be better than double then?
 - There is no data error in the operands
 - Alternatively, the errors cancel out themselves favorably

Overview

Why Bother with Mixed Precision?

Precision and Accuracy

Floating Point Operations

Mixed Precision Iterative Refinement

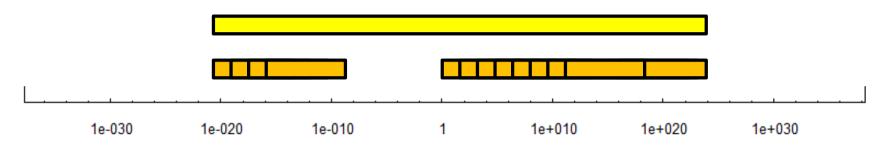


Understanding Floating Point Operations

- Number representation s23e8
 - $-a = | 1bit sign s_a | 23 bit mantissa m_a | 8 bit exponent e_a |$
- Multiplication a * b
 - Operations: $s_a^*s_b$, $m_a^*m_b$, $e_a^*+e_b^*$
 - Exact format: s46e9 = s23e8 * s23e8
 - Main error: Mantissa truncated from 46 bit to 23 bit
- Addition a + b
 - Operations: $e_{diff} = e_a e_b$, $m_a + (m_b >> e_{diff})$, normalize
 - **Exact format**: s278e8 = s23e8 + s23e8
 - Main error: Mantissa truncated from 278 bit to 23 bit

Commutative Summation

$$s = \sum_{i \in I} a_i$$



$$s = s_0 + s_1 + s_2$$

$$s_2 = \sum_{i \in I_2} a_i \qquad s_1 = \sum_{i \in I_1} a_i \qquad s_0 = \sum_{i \in I_0} a_i$$

$$1 = -030 \qquad 1 = -020 \qquad 1 \qquad 1 \qquad 1 = +020 \qquad 1 = +030$$

Commutative Summation Example

• 1 +
$$0.00000004 =_{db} 1.00000004 =_{fl} 1$$

In float s23e8

$$s = \Sigma a_i = \frac{1}{2} + \frac{1}{2} + 0.00000004 - 0.000000003 =_{fl} 1$$

In double s52e11

$$s = \Sigma a_i =_{db} 1.00000001$$

In mixed double/float

$$s_0 = \Sigma_0 a_i = \frac{1}{2} + \frac{1}{2} =_{fl} 1$$

 $s_1 = \Sigma_1 a_i = 0.00000004 - 0.00000003 =_{fl} 0.00000001$
 $s = s_0 + s_1 =_{db} 1.00000001$

Dependent Summation

$$s =_{db} s_0 + s_1 + s_2$$

$$s_2 =_{fl} \sum_{i \in I_2} a_i \qquad s_1 =_{fl} \sum_{i \in I_1} a_i \quad s_0 =_{fl} \sum_{i \in I_0} a_i$$

$$1 = -030 \qquad 1 = -020 \qquad 1 \qquad 1 \qquad 1 = +010 \qquad 1 = +020 \qquad 1 = +030$$

$$s =_{db} \sum_{i \in I} a_i$$
, $a_i = f_{db}(a_{i-1})$, with slow double precision $f_{db}()$

$$s_0 = \int_{I} \sum_{i \in I_0} a_i, \quad s_1 = \int_{I} \sum_{i \in I_1} a_i, \quad s_2 = \int_{I} \sum_{i \in I_2} a_i, \quad s = \int_{I} \sum_{i \in I_2} a_i, \quad s = \int_{I} \sum_{i \in I_2} a_i, \quad s = \int_{I} \sum_{i \in I} a_i, \quad s =$$

 $a_i = f_{fl}(a_{i-1})$, with fast single precision $f_{fl}()$

Overview

Why Bother with Mixed Precision?

Precision and Accuracy

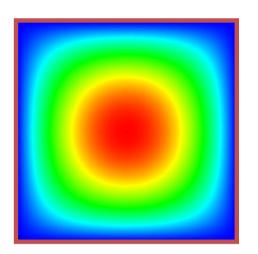
Floating Point Operations

Mixed Precision Iterative Refinement



PDE Example: Poisson Problem

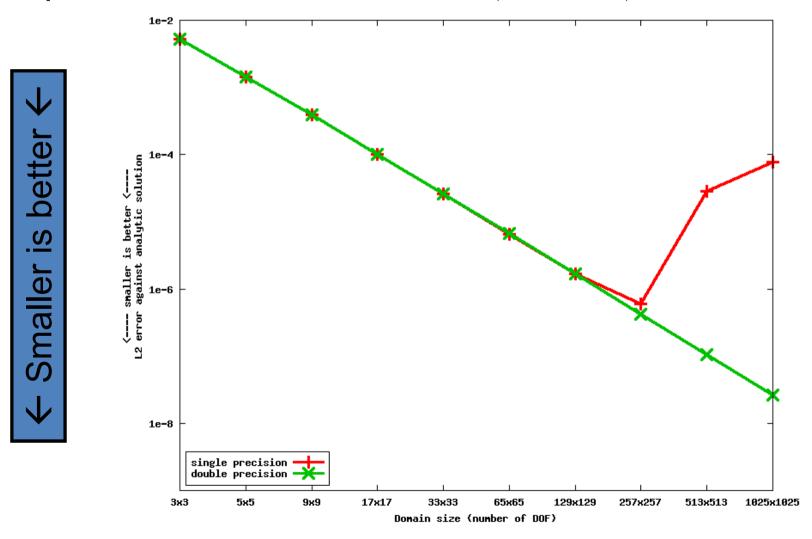
- - Δ u = f
- Unit square [0,1]^2
- Bilinear conforming FEs (Q1)
- Regular quadrilateral grid
- Zero Dirichlet BCs
- Analytic test function x(1-x)y(1-y)



Solved with multigrid until norms of residuals indicate convergence

PDE Example: Poisson Problem

- FEM theory: pure discretization error
- Expected error reduction of 4 (i.e. h^2) in each level



Mixed Precision Iterative Refinement

Exploit the speed of low precision and obtain a result of high accuracy

$$d_k = b-Ax_k$$

$$Ac_k = d_k$$

$$x_{k+1} = x_k + c_k$$

$$k = k+1$$

 $d_k = b-Ax_k$ | Compute in high precision (cheap) Solve in low precision (fast) $x_{k+1} = x_k + c_k$ Correct in high precision (cheap) k = k+1 Iterate until convergence in high Iterate until convergence in high precision

- Low precision solution is used as a preconditioner in a high precision iterative method
 - A is small and dense: Solve $Ac_k = d_k$ directly
 - A is large and sparse: Solve (approximately) $Ac_k = d_k$ with an iterative method itself

Direct Scheme for Small, Dense A

Algorithm

- Compute PA=LU once in single precision
- Use LU decomposition to solve $Ly=Pd_k$, $Uc_k=y$ in each step

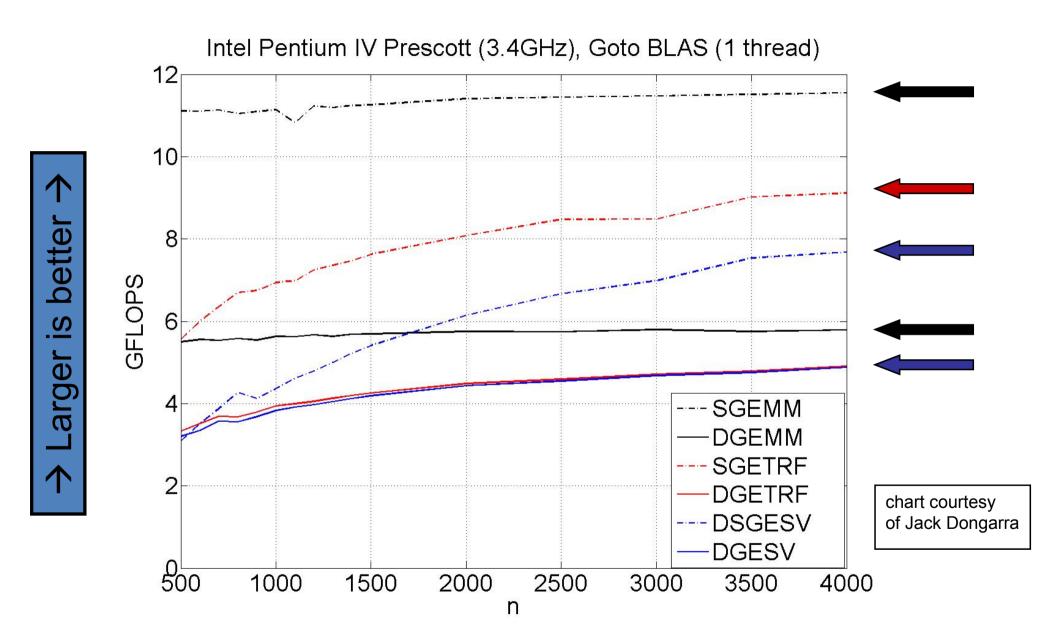
Main reasons for speedup

- Computation of LU decomposition is O(n^3)
- Computation of LU is much faster in single than in double
- Solution using LU for several RHS is only $O(n^2)$

Upper bound for iteration count

- $ceil(t_d/(t_s-K))$, where K,t_d,t_s are log10 of matrix condition and double and single precision (e.g. t_d approx 16)

CPU SSE Results: LU Solver



[Langou et al. Exploiting the performance of 32 bit floating point arithmetic in obtaining 64 bit accuracy (revisiting iterative refinement for linear systems), SC 2006]

Iterative Scheme for Large, Sparse A

Algorithm

- Inner solver: Conjugate Gradients, Multigrid
- Correction loop can run on CPU or on GPU (old GPUs: emulated precision; new GPUs: true double precision)
- Terminate inner solver after fixed number of iterations, fixed error reduction or convergence

Main reason for speedup

Inner solver on the GPU runs almost at peak bandwidth

Applicability

Works even for very ill-conditioned matrices

GPU Performance Results

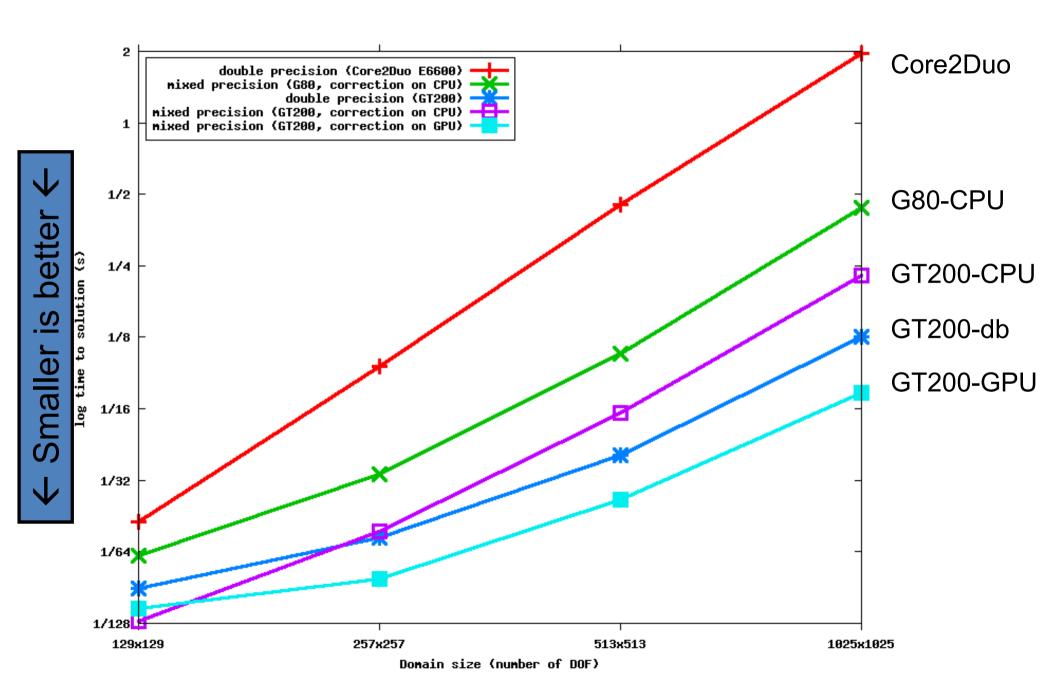
Test problem

- Poisson on unit square
- Multigrid solver
- N=33² to N=1025² DOF (=mesh points for Q1 FE)

Solver combination parameter space

- CPU implementation (Core2Duo E6600, SSE-optimized, double)
- CUDA implementation (GeForce 8800 GTX and GeForce GTX 280)
 - Mixed precision, correction on CPU (G80 and GT200)
 - Native double precision (GT200 only)
 - Mixed precision, correction on GPU (GT200 only)

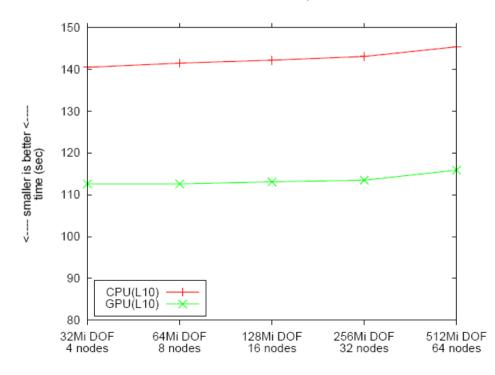
GPU Performance Results: CUDA

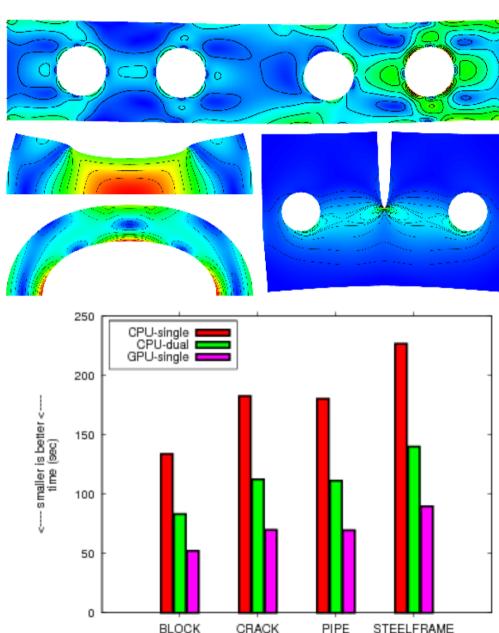


Mixed Precision on GPU Clusters

Total speedup of 2.7x for Quadro 5600 vs. AMD Santa Rosa (16 node cluster)

Good weak scalability on up to 64 nodes (dual Xeon, Quadro 1400 GPUs)





Conclusions

- The relation between computational precision and final accuracy is complicated but analyzable
- When single precision alone fails iterative refinement recovers the full accuracy with few double precision ops
- Mixed precision methods benefit bandwidth and even more computation bound algorithms
- Double precision GPUs are best utilized in mixed precision mode achieving outstanding performance and accuracy
- The benefits also extend to GPU clusters



Robert Strzodka, Max Planck Center, Saarbrücken, Germany Dominik Göddeke, Dortmund University of Technology, Germany

Mixed Precision Methods on GPUs

www.mpii.de/~strzodka/ www.mathematik.tu-dortmund.de/~goeddeke/

- D. Göddeke, R.Strzodka and S.Turek: *Performance and accuracy of hardware-oriented native-, emulated- and mixed-precision solvers in FEM simulations*, IJPEDS, 2007
- D. Göddeke and R.Strzodka, *Performance and accuracy of hardware-oriented native-, emulated- and mixed-precision solvers in FEM simulations (Part 2: Double Precision GPUs)*, TU Dortmund, Technical Report, 2008
- D. Göddeke, H. Wobker, R. Strzodka, J. Mohd-Yusof, P. McCormick, and S. Turek. *Co-processor acceleration of an unmodified parallel solid mechanics code with FEASTGPU*, IJCSE, 2008