

Accelerating the Multifrontal Method

Information Sciences Institute





22 June 2012 Bob Lucas, Gene Wagenbreth, Dan Davis, Roger Grimes {rflucas,genew,ddavis}@isi.edu and grimes@lstc.com



3D Finite Element Test Problem









LS-DYNA Implicit Half Million Elements



Timing information CPU(seconds) %CPU Clock(seconds) %Clock

Initialization .	2.6513E+01	0.24	2.7598E+01	1.53
Element proces	sing 8.1698E+0	1 0.73	7.6144E+0	1 4.22
Binary databas	ses 2.0471E+00	0.02	4.8657E-01	0.03
ASCII databas	e 1.3164E-04	0.00	1.5000E-05	0.00
Contact algorith	hm 2.9563E+01	0.26	5.7157E+00	0.32
Interf. ID	1 7.6828E+00	0.07	1.2124E+00	0.07
Interf. ID	2 4.8080E+00	0.04	9.7007E-01	0.05
Interf. ID	3 4.4934E+00	0.04	9.8310E-01	0.05
Interf. ID	4 5.5574E+00	0.05	1.209 <mark>1E+00</mark>	0.07
Interf. ID	5 7.0113E+00	0.06	1.33 <mark>80E+00</mark>	0.07
Contact entitie	es 0.0000E+00	0.00	0.0 <mark>000E+00</mark>	0.00
Rigid bodies	8.5901E-04	0.00	3. <mark>1600E-04</mark>	<mark>0.</mark> 00
Implicit Nonline	ear 8.7434E+00	0.08	6 <mark>.5756E+00</mark>	<mark>0.</mark> 36
Implicit Lin. Alg	g 1.1063E+04	98.67	1. <mark>6898E+03</mark>	<mark>93.5</mark> 5

Totals 1.1212E+04 100.00 1.8063E+03 100.00



Half million elements, eight Intel Nehalem cores

Major software components



- All of LS-DYNA
 - Linear Solver
 - Reordering (Metis) & symbolic factorization
 - Redistribute & permute sparse matrix
 - Factor, i.e., for all supernodes:
 - Assemble the frontal matrix
 - Factor the frontal matrix
 - Stack its Schur complement
 - Solve



USC Viterbi



Time in solver (sec) Eight Nehalem Threads



WCT: symbolic factorization =	18.952
WCT: matrix redistribution =	5.221
WCT: factorization = 795	.354
WCT: numeric solve = 2	2.705
WCT: total imfslv_mf2 = 82	26.735
WCT: symbolic factorization =	16.560
WCT: matrix redistribution =	5.087
WCT: factorization = 804	.471
WCT: numeric solve = 2	2.819
WCT total imfsly mf2 = 82	0 281







Start by accelerating factorization



- All of LS-DYNA
 - Multifrontal Linear Solver
 - Reordering (Metis) & symbolic factorization
 - Redistribute & permute sparse matrix
 - Factor, i.e., for all supernodes
 - Assemble the frontal matrix
 - Factor the frontal matrix with accelerator
 - Stack its Schur complement
 - Solve





Multifontal Method Toy Problem



X* X**XX

6







Multifrontal View of a Toy Matrix

USC Viterbi School of Engineering







Automotive Hood Inner Panel Springback using LS-DYNA



"Hood" Elimination Tree



Each frontal matrix's triangle scaled by operations required to factor it.



USC Viterbi

Exploiting Concurrency Toy with MPI & OpenMP







USC Viterbi



Why Accelerators?







Tesla vs Nehalem, courtesy of Stan Posey Similar expectations for Intel Phi

USC Viterbi School of Engineering Factorization of a Frontal Matrix



Assemble frontal matrix on host CPU

Initialize by sending panel of assembled frontal matrix

Only large frontal matrices due to high cost of sending data to and from accelerator







Eliminate panels



Factor diagonal block







Eliminate panels



Eliminate off-diagonal panel







Fill Upper Triangle











Update panels with DGEMM S = S - L * U

Assumes vendor math library DGEMM is extremely fast





USC Viterbi





Wider panels in Schur complement

USC Viterbi

School of Engineering

DGEMM is even faster S = S - L * U









Return error if diagonal of 0.0 encountered or pivot threshold exceeded

Otherwise complete frontal matrix is returned

Schur complement added to initial values on host CPU







Panels distributed with MPI

Factor panels broadcast and sent to GPU.

GPU eliminates factor panel and performs large-rank updates to the Schurcomplement.

Load imbalance is major bottleneck, so this needs revisiting





USC Viterbi



OpenMP on MIC (and on Xeon host)



```
do j = jl, jr
do i = jl + 1, ld
x = 0.0
do k = jl, j - 1
x = x + s(i, k) * s(k, j)
end do
s(i, j) = s(i, j) - x
end do
end do
```



```
sq = jr - jl + 1
c$omp parallel do
c$omp& shared (jl, jr, sq, ld, s)
c$omp& private (ii, j, i, il, ir)
    do ii = jl + 1, ld, sq
      do j = jl, jr
       \mathbf{i}\mathbf{l} = \mathbf{i}\mathbf{i}
       ir = min(ii + sq - 1, Id)
       do i = il, ir
        x = 0.0
         do k = jl, j - 1
          x = x + s(i, k) * s(k, j)
         end do
         s(i, j) = s(i, j) - x
       end do
      end do
    end do
```

Fortran vs CUDA



```
do j = jl, jr
do i = jr + 1, ld
x = 0.0
do k = jl, j - 1
x = x + s(i, k) * s(k, j)
end do
s(i, j) = s(i, j) - x
end do
end do
```

USC Viterbi



```
ip=0;
for (j = jl; j <= jr; j++) {</pre>
  if(ltid <= (j-1)-jl){
    gpulskj(ip+ltid) = s[IDXS(jl+ltid,j)];
  ip = ip + (j - 1) - jl + 1;
____syncthreads();
for (i = jr + 1 + tid; i <= ld;
     i += GPUL THREAD COUNT) {
  for (j = jl; j <= jr; j++) {</pre>
    gpuls(j-jl,ltid) = s[IDXS(i,j)];
    }
  ip=0;
  for (j = jl; j <= jr; j++) {</pre>
    x = 0.0f;
    for (k = jl; k <= (j-1); k++) {</pre>
      x = x + gpuls(k-jl, ltid) * gpulskj(ip);
      ip = ip + 1;
      gpuls(j-jl,ltid) -= x;
    }
  for (j = jl; j <= jr; j++) {</pre>
    s[IDXS(i,j)] = gpuls(j-jl,ltid);
    }
  }
```



Factorization Performance on Individual Frontal Matrices











Multithreaded Host Throughput Factoring a Model Frontal Matrix



Intel Nehalem



Host Throughput with MPI Factoring a Model Frontal Matrix

USC Viterbi











USC Viterbi School of Engineering Nehalem & Tesla Throughput Factoring a Model Frontal Matrix







For larger frontal matrices, >250 Gflop/s have been observed

USC Viterbi School of Engineering Factoring a Model Frontal Matrix



CPU THREAD, GPU, MIC PERFORMANCE





GPUs Get the Larger Supernodes







Relative Time (sec)



Benchmark Time in Seconds





500 CPU 1 Node CPU 2 Node GPU 1 Node GPU 2 Node 400 300 200 100 Û outer3 al16 outer2 **Nested cylinders benchmarks**

Benchmark GFLOPS



Work in Progress



- Near Term
 - Further accelerator performance optimization
 - Better MPI + accelerators
 - Pivoting on accelerators for numerical stability
- Longer Term
 - Factor many small frontal matrices on the accelerators
 - Assemble initial values
 - Maintain real stack
 - If the entire matrix fits on the accelerator
 - Forward and back solves



• Exploit GDRAM memory bandwidth



Research Initially Funded by US JFCOM and AFRL



This material is based on research sponsored by the U.S. Joint Forces Command via a contract with the Lockheed Martin Corporation and SimIS, Inc., and on research sponsored by the Air Force Research Laboratory under agreement numbers F30602-02-C-0213 and FA8750-05-2-0204. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the U.S. Government. Approved for public release; distribution is unlimited.





Bonus Slides







GPU Architecture



Courtesy NVIDIA

Multiple SIMD cores

USC Viterbi

School of Engineering

Multithreaded O(1000) per GPU

Banked shared memory 32 Kbytes Tesla 48 Kbytes Fermi

Simple thread model Only sync at host



A set of SIMD multiprocessors with on-chip shared memory.



Fortran vs CUDA



```
do j = jl, jr
do i = jr + 1, ld
x = 0.0
do k = jl, j - 1
x = x + s(i, k) * s(k, j)
end do
s(i, j) = s(i, j) - x
end do
end do
```

USC Viterbi



```
ip=0;
for (j = jl; j <= jr; j++) {</pre>
  if(ltid <= (j-1)-jl){
    gpulskj(ip+ltid) = s[IDXS(jl+ltid,j)];
  ip = ip + (j - 1) - jl + 1;
____syncthreads();
for (i = jr + 1 + tid; i <= ld;
     i += GPUL THREAD COUNT) {
  for (j = jl; j <= jr; j++) {</pre>
    gpuls(j-jl,ltid) = s[IDXS(i,j)];
    }
  ip=0;
  for (j = jl; j <= jr; j++) {</pre>
    x = 0.0f;
    for (k = jl; k <= (j-1); k++) {</pre>
      x = x + gpuls(k-jl, ltid) * gpulskj(ip);
      ip = ip + 1;
      gpuls(j-jl,ltid) -= x;
    }
  for (j = jl; j <= jr; j++) {</pre>
    s[IDXS(i,j)] = gpuls(j-jl,ltid);
    }
  }
```



USC Viterbi School of Engineering GPU gets larger Supernodes







6e+10 1 thread 2 threads 3 threads 4 threads 5e+10 5 threads 6 threads 7 threads 8 threads 4e+10 FLOPS per SEC 3e+10 2e+10 1e+10 Û 2 6 8 10 12 16 14 18 4

Multicore Performance per Level (CPU Only)



Level

Relative Performance (sec) half-million elements, in-core



MPI, OpenMP, GPU LS-DYNA 1st Factorization 8 threads 795 1806 8 threads + GPU 283 775 2 MPI x 4 threads 1659 772 2 MPI x 4 + 2 GPU 211 537



USC Viterbi

School of Engineering

Nested cylinders benchmark

USC Viterbi School of Engineering





MPI, OpenMP, GPU LS-DYNA 1st Factorization 8 threads 6131 15562 8 threads + GPU 2828 9143 2 MPI x 4 threads 15287 6248 2 MPI x 4 + 2 GPU 8214 2567



Nested cylinders benchmark

USC Viterbi School of Engineering Time in out-of-core solver (sec) Eight threads plus GPU

grep WCT mes0000 =>

WCT: symbolic factorization = 51.764
WCT: matrix redistribution = 10.875
WCT: factorization = 2828.933
WCT: numeric solve = 1409.523
WCT: total imfslv_mf2 = 4316.063
WCT: symbolic factorization = 40.244
WCT: matrix redistribution = 10.685
WCT: factorization $= 2679.422$
WCT: numeric solve = 1543.953
WCT: total imfslv_mf2 = 4346.260

Note: almost half the elapsed time for LS-DYNA (~4400 sec. out of 9143) is spent in disk I/O.

