Optimizing PLASMA Eigensolver on Large SGI UV Shared Memory Systems

Cheng Liao SR Software Engineer November, 2016



List of Eigensolvers

| SCALAPACK | ELPA | EIGENEXA | PLASMA |
|--|---|---|---|
| PDSYEVD, PDSYEVX, PDSYEV, PDSYEVR | ELPA1, EPLA2 | EIGEN_S, EIGEN_SX | DSYEVR (SMP only, no MPI, dataflow scheduler) |
| 1-stage tridiagonalization, 50% memory bound kernel | 1-stage & 2-stage, the latter tridiagonalizes much faster | improved 1-stage penta- diagonalization with sym2v | 2-stage tridiagonalization with tile storage layout |
| tridiagonal system solvers: D&C, BI, QR, MRRR | D&C tridiagonal solver | D&C pentadiagonal solver that takes longer to execute | LAPACK MRRR tridiagonal solver dstemr |
| compute bound 1-stage back-transformation | compute bound 1/2-stage back-transformation | compute bound 1-stage back-transformation | compute bound 2-stage back-transformation |
| de facto Industrial standards but slow | likes high FLOPS systems | likes high memory bandwidth systems | likes high FLOPS systems |

©2016 SGI

2

sgi

SGI UV3000 – Commodity CPUs + Proprietary Node Controllers with Routing Options for Cache Coherence



To backplane

To backplane

copy with shared memory, THP on

latency bet node 0 and node 0: 87.5ns latency bet node 0 and node 1: 400.5 ns latency bet node 0 and node 2: 507.8 ns latency bet node 0 and node 3: 507.7ns

bcopy on node 0 dest 0 src 0 BW: 6.407 GB/s bcopy on node 0 dest 0 src 1 BW: 2.465 GB/s bcopy on node 0 dest 0 src 2 BW: 1.977 GB/s bcopy on node 0 dest 0 src 3 BW: 1.981 GB/s

bcopy on node 0 dest 0 src 0 BW: 6.316 GB/s bcopy on node 0 dest 1 src 0 BW: 5.906 GB/s bcopy on node 0 dest 2 src 0 BW: 5.275 GB/s bcopy on node 0 dest 3 src 0 BW: 5.380 GB/s

3

52

...

©2016 SGI

Five Computational Phases of PLASMA Eigensolver



Sg



NUMA Considerations



 scratch buffers (MKL + PLASMA, etc.)

• A, V & T are in this category

Watch out for "OpenMP

Sg

ambush"

5

Arrange data and pin threads to minimize communication and avoid network hotspots!

©2016 SGI

Sample 'Best' NUMA Placement – Matrix Multiply on 4 Quad Core Nodes



memory placement color coded!

©2016 SGI

6

Sgi

Algorithmic Issue(1) – DBR Degree of Parallelism



HD-

HD.

©2016 SGI

Employing tree parallel QR with multiple Householder domains can improve dense to band reduction performance very significantly:



Number of HDs vs run time (secs)

Incremental Sequential QR Factorization As Illustrated in LAWN204



©2016 SG

8

Sgi

A Note on FP OP Counts

```
for (k = 0; k < NT-1; k++) // NT(==N/Nb): #tiles/row, BS: subdomain size
 BS = (NT - k - 2)/HD + 1; //HD: Number of Householder domains
//local QR factorizations on leading tiles of subdomains
 for (m = k+1; m < NT; m += BS) DGEQRT;
 // Apply the local reflectors
 // LEFT and RIGHT on the diagonal blocks
 for (m = k+1; m < NT; m += BS) DSYRFB;
 // RIGHT on tile column until the bottom
 for (m = k+1; m < NT; m += BS) { for (n = m+1; n < NT; n++) DORMQR; }
 // LEFT on tile row until the diagonal
 for (m = k+1; m < NT; m += BS) { for (n = k+1; n < m; n++) DORMQR; }
 // include other tiles in the subdomains
 for (M = k+1; M < NT; M += BS) {
  for (m =M+1; m < min(M+BS,NT); m++) {
    DTSQRT:
    for (i = k+1; i < m; i++) DTSMQR; // LEFT, excluding i=M
    for (j = m+1; j < NT; j++) DTSMQR; // RIGHT
    DTSMQRLR; // LEFT or RIGHT
 // tree-based merge of the local factors
 for (RD = BS; RD < NT-k-1; RD *= 2) {
  for (M = k+1; M+RD < NT; M += 2*RD)
    DTTORT:
    for (i=k+1; I<M+RD-1; i++) DTTMQR; // LEFT, excluding i=M
    for (j=M+RD+1; j <NT; j++) DTTMQR; // RIGHT
    DTTMQRLR; // LEFT or RIGHT
```

| Compute Kernel | Number of calls | | | | | | | |
|----------------|-----------------|----|--|--|--|--|--|--|
| dgeqrt | 1 | 2 | | | | | | |
| dsyrfb | 1 | 2 | | | | | | |
| dormqr | 5 | 10 | | | | | | |
| dtsqrt | 5 | 4 | | | | | | |
| dtsmqr | 20 | 16 | | | | | | |
| dtsmqrlr | 5 | 4 | | | | | | |
| dttqrt | n/a | 1 | | | | | | |
| dttmqr | n.a | 4 | | | | | | |
| dttmqrlr | n/a | 1 | | | | | | |
| | | | | | | | | |

Modest numbers of Householder domains do not change the complexity of D2B reduction:

$$\sum_{k=0}^{NT-2} \sum_{m=k+2}^{NT-1} \sum_{j=k+1}^{NT-1} 5Nb^3 \sim = \left(\frac{5}{3}\right) N^3$$

The 1st stage back-transformation complexity also does not change.

9

dormqr OPs need to be considered if the number of Householder domains is significant relative to number of tiles/row.

Algorithmic Issue(2)- Communication and Thread Placement in Bulge Chasing

| ٠ | | | | | | | | | | | 1 | E | 1 | | | E | Г | 1 | | | ٠ | | | | | | | | | | | • | 1 | | Γ | | | | | | | Г |
|---|----|----------|---|---|----------|---|---|---|---|---|----|---|---|---|---|---|-----------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ٠ | • | | | | | | | | | | Т | | | | | | \square | | | | ٠ | ٠ | | | | | | | | | | • | • | | | | | | | | | Г |
| ٠ | • | ٠ | | | | | | | | | | 6 | | | | | | | | | | ٠ | ٠ | | | | | | | | | | • | • | | | | | | | | Г |
| ٠ | • | ٠ | • | | | | | | | | • | • | 6 | | | | | | | | | ٠ | ٠ | ٠ | | | | | | | | | • | ٠ | • | | | | | | | Г |
| | • | • | • | • | | Г | | | | | Ie | • | • | ٠ | | | Г | | | | | ٠ | • | • | 6 | | | | | | | Г | • | • | • | • | | | | | | Г |
| | | ٠ | • | • | ٠ | | | | | | • | • | • | ٠ | • | | \square | | | | | • | ٠ | • | • | 0 | | | | | | | T | ٠ | • | • | • | | | | | Г |
| | | | • | • | ٠ | • | | | | | • | • | • | ٠ | • | • | | | | | | • | ٠ | • | • | • | 0 | | | | | | Г | ٠ | • | • | • | • | | | | Г |
| | | Г | Г | • | ٠ | • | • | Г | | | | | | • | • | • | • | | | | | | | | • | • | • | • | | | | Г | Г | | | • | • | • | 2 | | | Г |
| | | | | | • | • | • | • | | | | | | | • | • | • | ٠ | | | | | | | ٠ | ٠ | ٠ | • | ٠ | | | | | | | • | • | • | ٠ | 0 | | Γ |
| | | | | | | • | • | • | • | | | | | | | • | • | • | • | | | | | | ٠ | ٠ | ٠ | • | ٠ | • | | | | | | • | • | • | ٠ | ٠ | 0 | Γ |
| | L. | [| Г | Γ | — | Г | • | • | • | ٠ | T | | | | | | • | ٠ | • | • | | | | | | | | • | • | • | • | Г | Т | | Γ | | | | ٠ | ٠ | ٠ | |

PLASMA implementation:

- All work is done by 3 types of compute kernels/tasks.
- T?y(z-1) and T?(y-1)(z+2) done before T?yz.

Shortcomings:

- Left and right applications of Householder reflectors need to load matrix blocks 4 times.
- No thread placement strategy.

Improved communication:

$$(\mathbf{I} - \dot{\alpha} \bullet \boldsymbol{u} \bullet \boldsymbol{u}^{t}) \bullet \mathbf{B} \bullet (\mathbf{I} - \tau \bullet \mathbf{v} \bullet \mathbf{v}^{t})$$

(4x loads of **B**)

B - $\dot{\alpha} \bullet u \bullet u^{t} \bullet B + WORK \bullet v^{t}$

Where $w = B \bullet v$, $b = u^t \bullet w$, $s = \dot{\alpha} \bullet \tau \bullet b$, and WORK = $-\tau \bullet w + s \bullet u$. (column-wise, 2x loads of B) Improved performance:

| Problem/ | Bulge | Bulge Chasing performance on 60 e5-4627 v3 sockets | | | | | | | | | | | | | |
|---------------------|--------|--|--------------|---------------------|----------------------------------|--------|--|--|--|--|--|--|--|--|--|
| Tile Size | Origir | nal code | Improv th | /ed+linear reads | Improved+round- robin threads | | | | | | | | | | |
| | sec | gflops | sec | gflops | second | gflops | | | | | | | | | |
| N=115200, NB=480 | 106 | 360 | 74 | 518 | 78 | 490 | | | | | | | | | |
| N=115200,N B=640 | 173 | 295 | 139 | 366 | 109 | 468 | | | | | | | | | |
| N=115200,N B=960 | 684 | 112 | 458 | 167 | 240 | 318 | | | | | | | | | |
| N=288000,N B=960 | 7676 | 62 | 8917 | 54 | 1387 | 345 | | | | | | | | | |

Round-Robin Thread Placement:



number, *y* represents the bulge chasing sweep number and *z* is the taskid of the sweep.

| Socket #0 | Socket #1 | Socket #2 | Socket #3 |
|-----------|-----------|-----------|-----------|
| Thread 0 | Thread 1 | Thread 2 | Thread 3 |
| Thread 4 | Thread 5 | Thread 6 | Thread 7 |
| | | | |

©2016 SGI

Algorithmic Issue(3)- 1D Parallel Decomposition in Eigenvector Back-transformations

It is straightforward to implement 1D parallel decomposition for eigenvector backtransformation. However it also is well known that 1D parallel decomposition has a high communication to computation ratio. The problem size needs be sufficient large to ensure good performance.

$$\sum_{j=vblksiz}^{N,vblksiz} \frac{j}{NB} (4N(NB + vblksiz)vblksiz + Nvblksiz^2)$$

$$\sim = 2N^3 \left(1 + 1.25 \frac{vblksiz}{NB}\right)$$

Per Socket Communication for Q2 BT:

11



©2016 SGI

 $\frac{N^{2}}{2} * size of (double) + \frac{N^{2} v b l k s i z}{2 * N B} * size of (double) = \sim 5N^{2} b y tes$ **V T**

Communication/computation ratio is roughly $2/\beta$: $\frac{N}{\alpha}_{NP}$, where α , β , NP are the BW, FP speed and number of sockets

PLASMA Eigenvector Back-transformations



Illustrations from Haidar, Luszczek and Dongarra 2014 paper

Sgl



Effects of Tile & Memory Page Sizes



©2016 SGI

13

Sgl

Performance Comparisons



Elapsed times of N=64000 on 60 sockets



Elapsed times of N=115200 on 60 sockets



Elapsed times of N=288000 on 60 sockets



PLASMA GFLOPS improvement w.r.t. problem size

©2016 SGI

Conclusion

The PLASMA Eigensolver works well with large problem sizes. Can the strengths of PLASMA and ELPA be combined?

15

Sg