Petascale in CFD Cetraro, 29th of June, 2011

Christian Simmendinger, T-Systems Solutions for Research



Overview

The TAU Solver



Scalability

- Core (Instruction Level)
 - SIMD SRC-SRC
- Thread (Task Level)
 - Beyond OpenMP
- System Level
 - GASPI / PGAS API

Maintainability of the CodeProgrammabilityPerformance Portability





The TAU Solver Vision: Digital Aircraft

TAU: 3D Finite Volume Reynolds Averaged Navier Stokes (RANS) Solver





SSE2/SSE4.2/AVX/MIC: From 128 bit SIMD to 512 bit SIMD. 16 DP Flop/Cycle – or just one.

Statistical Fact:

Compiler performance doubles every 18 years. Do not expect compilers to resolve the SIMD issue for you. Especially C/C++ takes a severe performance hit due to aliasing.



MaintainabilityProgrammabilityPerformance-Portability

Scalability

- Core (Instruction Level)
 - SIMD SRC-SRC
- Thread (Task Level)
 - Beyond OpenMP
- System Level
 - GASPI / PGAS API



Programmability , Maintainability and Performance Portability: Instruction Level

שב דטאבטירפערנאט-2חס.כ - בכנוףגפ אטא	TRUCTURED, 2nd.scout.c - Eclipse SDK
Search Project Run Window Help	ate Search Proj Window Help
erver 🗞 🛪 🗞 ד 🖓 ד 🖓 🖉 🗊 🖉 😂 🔗 ד 🖉 ד 🌤 🐤	rt server 🗞 🛪 🗞 🗸 🖓 🔻 💋 🖉 🗐 🗊 🖉 😂 🔗 🛪 🖓 🛪 🏷 🗇 🛪
🗈 traceControl.h 🛛 🖻 feufl3D-2nd.scout.c	🗈 finit.c 🕞 trace 📄 Makefile.sse4 🔒 feufl3D-2nd.scout.c 🛛 📄 M
<pre>Fint layerSide = SIDEJ_KDIM * SIDEJ_JDIM; FphysFlow uL, uR; FavVal roe; Fvector a, charVal, PhiVar; Fmatrix R; #pragma scout loop vectorize for (k = sk; k < ek; k++) { Fint indexNodeJp1 = l * layerNode + (j) * kDimNode + k; Fint indexNodeJp1 = indexNodeJp1 - kDimNode; Fint indexNodeJp2 = indexNodeJp1 + kDimNode; Fint indexNodeJp2 = indexNodeJp1 + kDimNode; Fint indexNodeJp2 = indexNodeJp1 + kDimNode; Fint indexSide = l * layerSide + j * kDimSide + k; Ffloat omegaRSqr = sqr(OMEGA * P_SIDE_R_J[indexSide].radius); Ffloat inv_A_norm = 1. / P_SIDE_A_J[indexSide].A_norm; musclVanLeer3D(data, &uL, &uR, omegaRSqr, indexNodeJm1, indexNode RoesAverage3DRel(data, &uL, &uR, &roe, omegaRSqr, &BETA_PRECOND); eigenvalue3D(&a, &P_SIDE_A_J[indexSide], &roe); alpha3D(data, &charVal, &uL, &uR, &roe, &P_SIDE_A_J[indexSide]); RPulliam3D(data, &R, &roe, P_SIDE_A_J[indexSide].A.x * inv_A_norm P_SIDE_A_J[indexSide].A.z * inv_A_norm); PhiRoe3D(&PhiVar, &a, &charVal); FluxPHalfRoe3D(&R, &PhiVar, &P_SIDE_A_J[indexSide], &P_SIDE_FLUX_} } </pre>	<pre>art_vect372 art_vect683 =</pre>
VOIA	





Mutual Exclusion

for (color = fcolor; color != NULL; color = color->succ)
 for (face = color->start; face < color->stop; face++)



```
for(color=get_color();color != NULL; color=get_next_color(color))
    for (face = color->start; face < color->stop; face++)
```

Maintainability

- Programmability
- Performance Portability

OpenMP - The Fork Join Model

- Incremental Parallelisation
- Amdahl Trap (Barrier and Load Imbalance)







Mutual Completion

```
for(color=get_color();color != NULL; color=get_next_color(color))
    for (face = color->start; face < color->stop; face++)
BARRIER;
for(pnt = 0; pnt < nown; pnt++)</pre>
```

♣

```
for(color=get_color();color != NULL; color=get_next_color(color))
    for (face = color->start; face < color->stop; face++)
while(get_pnt_range(&pstart, &pstop))
    for (pnt = pstart; pnt < pstop; pnt++)...</pre>
```

Maintainability

- Programmability
- Performance Portability

This Programming Model is Generic.

for (k = 0; k < KDIM; k++)
for (j = 0; j < JDIM; j++)
for (i = 0; i < IDIM; i++) {};
BARRIER;
for (k = 0; k < KDIM; k++) ...</pre>



```
for(k = get_k_index(); k != -1; k = get_next_k_index(k))
for (j = 0; j < JDIM; j++)
for (i = 0; i < IDIM; i++) {};
for(k = get_k_index(); k != -1; k = get_next_k_index(k)) ...</pre>
```

J.Jägersküpper (DLR), C.Simmendinger (T-Systems SfR). A Novel Shared-Memory Thread-Pool Implementation for Hybrid Parallel CFD Solvers, EuroPar 2011, (to appear).

This extension of the OpenMP Fork-Join Programming Model is

- fully asynchronous
- data flow driven
- implicitly load balanced
- exclusively uses local locks instead of global barriers
- expected to scale to O(10^2) cores





In a Partitioned Global Address Space every thread can read/write the entire global memory of the application.

 Opportunities for improved programmability (PGAS languages like UPC, CAF, Chapel and others)
 Opportunities for improved Scalability (GPI (Fraunhofer ITWM), GASPI)

Global Address Space Programming Interface (GASPI)

A PGAS API which aims at replacing bulk-synchronous communication with asynchronous, 1-sided, RDMA driven communication patterns.

Scalability

- Core (Instruction Level)
 - SIMD SRC-SRC
- Thread (Task Level)
 - Beyond OpenMP
- System Level
 - GASPI / PGAS API



Overlapping Communication and Computation

```
if (this_thread_exchanges_pnt_data(g)) /* lock halo */
{
    exchange_pointdata();
    post_exch_pnt_data(); /* unlock halo*/
}
```

Programmability, Maintainability and Performance Portability: System



Multithreaded send/recv

if ((sid = get_seq_num_and_lock(g)) < num_send_threads(g))
 send_dbl_threaded(comap, comap->ncommdomains, comap->commpartner,
 data, dim2, key);

If ((sid = thread_seq_num()) < num_recv_threads(g))
 recv_dbl_threaded(g, num_recv_threads, sid,
 comap, comap->ncommdomains, comap->commpartner, data, dim2, key);

Strong Scaling Use Case: How far can we scale?

4W Multigrid, 1 Iteration Step.

• 1 Allreduce

• 140 Next Neighbour Halo Exchanges









Summary



Scalability

- Core (Instruction Level)
 - SIMD SRC-SRC
- Thread (Task Level)
 - Beyond OpenMP
- System Level
 - GASPI / PGAS API / GASPI



We are confidemt, that we can reach O(Petascale) within the constraints of a typical production run.

