# CoreTSAR: Adaptive Worksharing for Heterogeneous Systems

Tom Scogland, Wu-chun Feng Dept. of Computer Science, Virginia Tech Barry Rountree, Bronis R. de Supinski Lawrence Livermore National Laboratory

Edited and presented by Harsh Khetawat





# Heterogeneity in HPC: Accelerators on the Top500



- Heterogeneity becoming the norm in HPC
- Top 10 HPC systems dominated by accelerated systems

### How do we target them?





# General Matrix Multiplication(GEMM): CPU serial

```
void runGemm(T *a, T *b, T *c) {
 for (int i = 0; i < N; i + +) {
  for (int j=0; j < N; +j) {
  c[(i * N) + j] *= B;
   for (int k = 0; k < N; ++k) 
    c[(i*N)+j] = A * a[(i*N)+k] * b[(k*N)+j];
                                   Targets: One CPU core
```







# **GEMM:** OpenMP

```
void runGemm(T *a, T *b, T *c) {
#pragma omp parallel for
for (int i = 0; i < N; i + ) {
 for (int j = 0; j < N; ++j) {
  c[(i * N) + j] *= B;
   for (int k=0; k < N; ++k) {
    c[(i*N)+j] = A * a[(i*N)+k] * b[(k*N)+j];
                                 Targets: All local CPU cores
```











Worker threads





# **GEMM: CUDA** (minimal)

```
global void cudag(T*a, T*b, T*c, TA, TB, int n) {
uint i = blockIdx.x * blockDim.x + threadIdx.x;
if(i < n)
 for (int j=0; j < n; ++j) {
  c[(i^*N)+j] = B;
  for (int k=0; k < n; ++k) {
   c[(i*N)+j] = A*a[(i*N)+k]*b[(k*N)+j];
} } } }
void runGemm(T **a, T **b, T **c) {
T^*ca, *cb, *cc; dim3 dB, dG;
size t size = N*N*sizeof(T);
dB.x=64; dB.y=dB.z=1;
dG.x=(N/dB.x)+1; dG.y=dG.z=1;
cudaMalloc(&ca, size);
cudaMalloc(&cb, size);
cudaMalloc(&cc, size);
cudaMemcpy(ca,*a,size,cudaMemcpyHostToDevice);
cudaMemcpy(cb,*b,size,cudaMemcpyHostToDevice);
cudaMemcpy(cc,*c,size,cudaMemcpyHostToDevice);
cudag \ll dG, dB \gg (a, b, c, A, B, N);
cudaMemcpy(*c,cc,size,cudaMemcpyDeviceToHost);
```

*lirginia*Tech

Invent the Future

#### Targets: **One** GPU



# **CUDA** Threading Behavior



# GEMM: Accelerated OpenMP (OpenMP 4.0 syntax)

```
void runGemm(T *a, T *b, T *c) {
#pragma omp target teams distribute parallel for
           map(tofrom: c[0:N][0:N]) \
    map(to: a[0:N][0:N], b[0:N][0:N])
for (int i = 0; i < N; i + +) {
  for (int j = 0; j < N; ++j) 
  c[(i^*N)+j] *=B;
   for (int k = 0; k < N; ++k) 
    c[(i*N)+j] = A * a[(i*N)+k] * b[(k*N)+j];
                                Targets: One accelerator or
                                       one CPU core
```





# Accelerated OpenMP Threading Behavior

#pragma acc kernels for...



1872

#### Accelerated OpenMP + OpenMP Threading Behavior



#### The Goal:

Work-share a Target Region Across the Whole System



#### Issues

- Computational Overhead
  - Launching tasks onto GPUs is **expensive** (compared to CPUs)
- Heterogeneity of Computational Devices
  - Mapping the right amount of work to the right core
- Memory Incoherence
  - CPUs and GPUs generally do *not* share memory





# Solutions

- Adaptive scheduling to minimize launch overhead (Splitter)
- Performance-prediction for load-balancing(CoreTSAR)
- Multi-device memory management (CoreTSAR)





# Background: GEMM: Our Previous Approach, Splitter

```
void runGemm(T *a, T *b, T *c) {
splitter * s = split init(N, SPLIT DYNAMIC, NULL, NULL);
for (int d it=0; d it<s->d end; d it++) {
 s = split next(no, d it);
#pragma omp parallel num threads(2)
  int tid = omp get thread num();
   if(tid == 0) split gpu start(s); else split cpu start(s);
  int start = tid == 0?s->gts : s->cts;
  int end = tid = 0?s->gts : s->cts;
#pragma omp target teams distribute parallel for
  map(tofrom: c[start:end][start:end]) \
  map(to: a[startend][startend])
  map(to: b[start:end][start:end])
  for (int i = start; i < end; i++) {
   for (int j=0; j < N; ++j) {
    c[(i^*N)+j] = B;
    for (int k=0; k< N; ++k) {
     c[(i*N)+j] = A*a[(i*N)+k]*b[(k*N)+j];
   } } }
  if(tid == 0) split gpu end(s); else split cpu end(s);
} } }
```

**VirginiaTech** 

Targets: **One** local accelerator **and** all CPU cores



# Background: Splitter Threading Behavior



### Background: Splitter's Main Contribution: Adaptive Static Scheduling Policies

- Minimize blocking time
  - No device should stand idle waiting for others
- Avoid overhead of running more tasks than necessary
  - Launching a task on a GPU, and moving memory, is expensive





# Background: Scheduling Policies: Static

- Divide work based on a static ratio
  - Straightforward extension of OpenMP static to divide work unevenly among targets, based on floating point performance by default
  - Pros: No scheduling overhead

- Cons: Cannot adapt if the ratio is wrong or situation changes Initial Ratio: 0.75



# Background: Scheduling Policies: Adaptive

- Generate a ratio based on behavior each pass
  - Predicts the runtime on each device and minimizes blocking time
  - Pros: Creates one task per device, can adapt to change
  - Cons: Only adapts at entry to the region



# Background: Scheduling Policies: Split

- Breaks each pass into sub-passes and adapts for each
  - Takes an extra parameter, "div," for number of sub-passes to create
  - Pros: Adapts faster and more often
  - Cons: Higher overhead





# Background: Scheduling Policies: Quick

- Breaks the **first** pass into sub-passes to train
  - Uses a short sub-pass in the first pass to train, then uses adaptive
  - Pros: Adapts earlier, keeps overhead similar to adaptive
  - Cons: May be mis-trained by the small sub-pass



# Background: Scheduling Policies: Quick

• Breaks the first pass into sub-passes to train

Invent the Future

Uses a short sub Pros. **Assumes each iteration** Initia does similar computation across 500 passes! Original/Master thread Worker threads Parallel region Accelerated region Reschedule

synergy.cs.vt.edu

# **GEMM: Accelerated OpenMP**

```
void runGemm(T *a, T *b, T *c) {
#pragma omp target teams distribute parallel for \
           map(tofrom: c[0:N][0:N]) \
    map(to: a[0:N][0:N], b[0:N][0:N])
for (int i = 0; i < N; i + +) {
  for (int j = 0; j < N; ++j) {
  c[(i^*N)+j] *=B;
   for (int k = 0; k < N; ++k) 
    c[(i*N)+j] = A * a[(i*N)+k] * b[(k*N)+j];
```





# GEMM: CoreTSAR Extended Accelerated OpenMP

```
void runGemm(T *a, T *b, T *c) {
\#pragma omp target teams distribute parallel for \setminus
    map(partial, tofrom: c[true:N][false:N]) \
    map(partial, to: a[true:N][false:N])
    map(to: b[0:N][0:N])
    hetero(true, all, adaptive)
 for (int i = 0; i < N; i + +) {
  for (int j = 0; j < N; ++j) {
   c[(i^*N)+j] = B;
   for (int k = 0; k < N; ++k) 
    c[(i*N)+j] = A * a[(i*N)+k] * b[(k*N)+j];
```

Targets: All local accelerators and CPU cores





#### **CoreTSAR** Threading Behavior



# Solutions

- Adaptive scheduling to minimize launch overhead (Splitter)
- Performance-prediction for load-balancing(CoreTSAR)
- Multi-device memory management (CoreTSAR)





# Performance Prediction: Our Previous Approach

- Designed to "split" work between CPUs and one GPU
- Used a simple extrapolation to balance predicted runtime







# Performance Prediction: Take Two

- Integer optimization program
  - globally optimal distribution based on the prediction
  - Minimize difference in runtime 1 dovices

 $I = \text{tot} \qquad \text{Issue: It is slow} \\ i_j = \text{itera} \\ f_j = \text{fraction or } \\ p_j = \text{recent time/iteration for device j} \qquad i_2 * p_2 - i_1 * p_1 = t_1^+ - t_1^- \\ n = \text{number of devices} \\ i_3 * p_3 - i_1 * p_1 = t_2^+ - t_2^- \\ \end{cases}$ 

$$t_i^+$$
(or  $t_i^-$ ) = time over (or under) equal

Invent the Future

 $i_n * p_n - i_1 * p_1 = t_{n-1}^+ - t_{n-1}^-$ 

n-1



### Integer Linear Program Performance



# What's Taking so Long?

# For a solve with 8 devices



# Performance Prediction: Take Three

• Linear optimization program

nvent the Future

- Solves for a **percentage** of iterations rather than total number
- Guaranteed to be within I = tot I

$$f_n * p_n - f_1 * p_1 = t_{n-1}^+ - t_{n-1}^-$$



### Linear Program Performance



# Solutions

- Adaptive scheduling to minimize launch overhead (Splitter)
- Performance-prediction for load-balancing(CoreTSAR)
- Multi-device memory management (CoreTSAR)





# Memory Association:

- Accelerated OpenMP:
  - Maps a logically-contiguous block of memory to a compute region
  - Can target only one device per thread at a time
- The alternative (or extension)
  - Specify association between iterations and data
  - Compute input/output set from assigned iterations in the runtime
  - Reduce memory transfers by only copying assigned data
  - Single directive is sufficient for any number of devices
  - Maintain consistency by merging **all** output into host memory at the end of each region





# Memory Association Syntax

#### map(partial, <direction>: <variable>[<assoc.>:<length>:<boundary>])

- <variable>: An array, matrix or pointer that conforms to the requirements of the copy, copyin, or copyout clauses
- <assoc.>: Whether to associate the dimension with the accelerator
- <length>: The number of elements in the dimension
- <boundary>: Number of "boundary" elements required in this dimension





# Memory Association Example: GEMM







### Memory Association Example: What Actually Happens?



synergy.cs.vt.edu



# Memory Association Example: Scheduling



### Memory Association Example: Data Distribution: Input



### Memory Association Example: Data Distribution: Output



# CoreTSAR Memory-Association Example: Simple Column-Wise Association

Directive: #pragma acc region hetero(TRUE) \ pcopy(mat[false:10][true:10])

Memory not used on this device I Input only Output only

Invent the Future



41

Synergy.cs.vt.edu

Input and output

# CoreTSAR Memory-Association Example: Row-Wise Single-Boundary

Directive: #pragma acc region hetero(TRUE) \ pcopy(mat[true:10:1][false:10])



### Results: Benchmarks Representative Subset

- GEMM Few passes
  - Matrix multiplication benchmark from PolyBench/GPU
- K-Means Few passes
  - Iterative clustering of points
- Helmholtz-Few passes, GPU Unsuitable
  - Jacobi iterative method implementing the Helmholtz equation
- CORR Few passes, heterogeneous iterations
  - Upper-triangular correlation matrix solver from PolyBench/GPU, all iterations do different amounts of work





# **Results: Experimental Setup**

- System
  - 2x 4-core Intel X5550 CPUs
  - 4x NVIDIA Tesla C2070 GPUs
  - Debian Squeeze Linux
  - PGI Accelerator compiler version 12.9
- Procedures
  - All parameters default, unless otherwise specified
  - Results represent 5 or more runs





#### Results for Co-Scheduling Amenable Benchmarks: GEMM: Dense Linear Algebra, Matrix Multiplication



#### Results for Co-Scheduling Amenable Benchmarks: GEMM: Dense Linear Algebra, Matrix Multiplication



# Results for Co-Scheduling Amenable Benchmarks: K-Means: Clustering Algorithm



#### Results for Co-Scheduling Averse Benchmarks: Helmholtz: Jacobi Solver of the Helmholtz Equation



#### Results for Co-Scheduling Averse Benchmarks: Helmholtz: With and Without GPU Back-off

With GPU-backoff

Implementation / Original

VirginiaTech

Invent the Future

1872



Synergy.cs.vt.edu

### **Results for Co-Scheduling Averse Benchmarks:** PolyBench CORR: Upper Triangular Matrix, Correlation





Every iteration, and range, does a different amount of work!







# CORR Scheduling Behavior with Adaptive

Uneven work between iterations causes oscillations in Adaptive

CPU

GPU

Device type

Invent the Future



SyNeRG synergy.cs.vt.edu

# **CORR Scheduling Behavior with Split**

#### Split does better, by not using the CPU cores!

CPU

GPU

Device type

Invent the Future





CORR with Adaptive GPU-only





# Comparison with State of the Art

- Chose two popular task schedulers
  - OmpSs
  - StarPU
- Both are general *task* schedulers
  - Arbitrary graphs of dependent tasks can be expressed
  - Scheduling is focused on distributing discrete *tasks* rather than loop iterations
- Our approach is complementary to theirs





# Comparison with State of the Art: Setup

- Ported three of the benchmarks to the OmpSs and StarPU task schedulers
  - CUDA/c versions of each kernel function added
  - OmpSs: Versioning stack scheduler used to allow alternatives, automatic scheduling between CPUs and GPUs
  - StarPU: directly implemented in c API using the history based performance model, primed with 10 runs, and the dmda scheduler
- Created CoreTSAR scheduled versions using the CUDA/C kernels
  - CoreTSAR does *not* require accelerated OpenMP regions
  - Allows direct 1-1 comparison with identical compute regions
- All compiled with –O3 on:
  - nvcc: CoreTSAR and StarPU
  - mnvcxx: OmpSs







# Conclusions

- Adaptive scheduling and dynamic task granularity provided speedups of as much as 2x over alternative methods
- Our dynamic schedulers achieve high performance across various applications and highly heterogeneous machines
- Heterogeneous task scheduling for accelerated OpenMP is a feasible and useful extension

# Questions?





# Future Directions

- Explore other accelerators such as FPGAs
  - OpenCL vs Synthesis
- Cost of memory movement not included in model
  - Stencil codes need boundaries which are calculated in each iteration
  - Cost vs benefit of moving boundaries between accelerators in each iteration
  - Asynchronous coherence?
- Uneven work in each iteration across passes
  - Application profiling?



