# A New, Portable Algorithm Framework for Parallel Linear Recurrence Problems

Wen-mei W. Hwu

with
Tom Jablin, Liwen Chang, Chris Rodrigues,
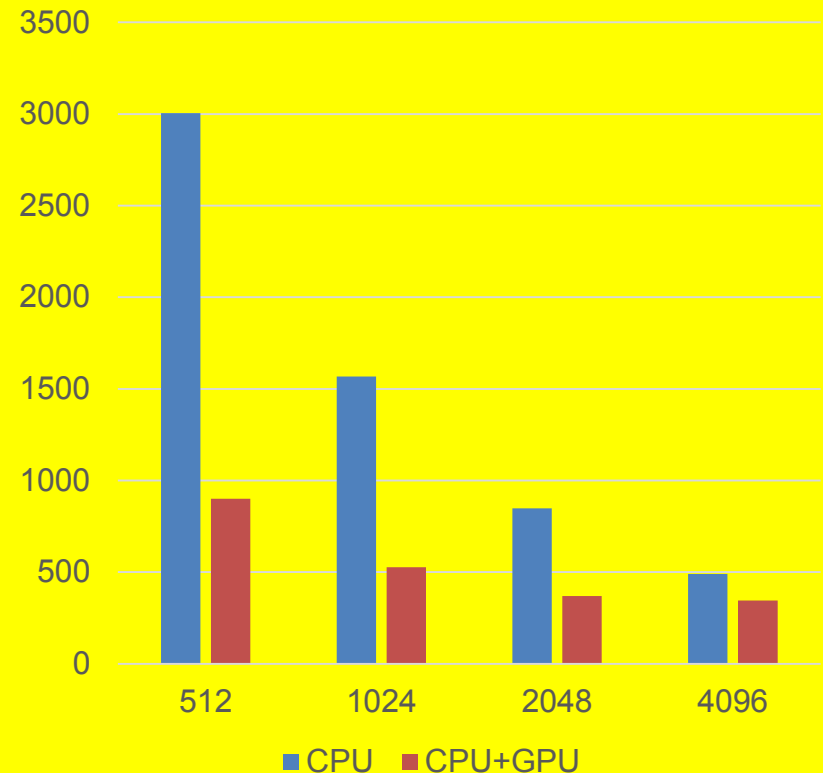Steven ShengZhou Wu, Abdul Dakkak

University of Illinois at Urbana-Champaign

# Two Current Challenges for Petascale GPU Computing

- At scale use of GPUs
  - Communication costs dominate beyond 2048 nodes
  - E.g., NAMD Limited by PME
  - Insufficient computation work
- Programming Efforts
  - This talk

Blue Waters K7 Nodes NAMD Strong Scaling – 100M Atoms



ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN
1867

SC13

# Writing efficient parallel code is complicated.
## Tools can provide focused help
## or broad help

**Planning how to execute an algorithm**   **Implementing the plan**

• Choose data structures

• Memory allocation
• Data movement
  GMAC

• Pointer operations
• Index arithmetic
  DL

Triolet, X10, Chappel, Nesl, DeLite, Par4All

• Map work/data to tasks
• Schedule tasks to threads

• Kernel dimensions
• Thread ID arithmetic
  OpenACC/ C++AMP/ Thrust

• Synchronization
• Temporary data structures

Tangram

Joint Lab 2013

# Levels of GPU Programming Languages

**Prototype & in development**                                   X10, Chapel, Nesl, Delite, Par4all, Triolet...

Implementation manages GPU threading and synchronization invisibly to user

**Next generation**                        OpenACC, C++AMP, Thrust, Bolt

Simplifies data movement, kernel details and kernel launch

Same GPU execution model (but less boilerplate)

**Current generation**                        CUDA, OpenCL, DirectCompute

# Where should the smarts be for Parallelization and Optimization?

- General-purpose language + parallelizing compiler
  - Requires a very intelligent compiler
  - Limited success outside of regular, static array algorithms

- Domain-specific language + domain-specific compiler
  - Simplify compiler's job with language restrictions and extensions
  - Requires customizing a compiler for each domain

- Parallel meta-library + general-purpose compiler
  - Library embodies parallelization policies and decisions
  - Uses a general-purpose compiler infrastructure
  - Extensible—just add library functions
  - Historically, library is the area with the most success in parallel computing

# Triolet – Composable Library-Driven Parallelization

- Allows library to collect multiple parallel operations and create an optimized arrangement

  - **Lazy evaluation** and aggressive inlining

  - **Loop fusion** to reduce communication and memory traffic

  - **Array partitioning** to reduce communication overhead

  - Library source-guided **parallelism optimization** of sequential, shared-memory, and/or distributed algorithms

- Loop-building decisions use information that is often known at compile time

  - By adding static typing to Python

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN
1867

# Example: Correlation Code

Compute f(x,y) for every x in xs and for every y in ys (Doubly nested loop)

```
def correlation(xs, ys):
    scores = (f(x,y) for x in xs for y in ys)
    return histogram(100, par(scores))
```

Compute it in parallel

Put scores into a 100-element histogram

# Triolet Compiler Intermediate Representation

- List comprehension and `par` build a package containing
    1. Desired parallelism
    2. Input data structures
    3. Loop body
    
    for each loop level
- Loop structure and parallelism annotations are **statically known**

```
correlation xs ys =
  let i = IdxNest HintPar
              (arraySlice xs)
              (λx. IdxFlat HintSeq
                        (arraySlice ys)
                        (λy. f x y ) )
  in histogram 100 i
```

Outer loop

Inner loop

Body

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN
1867
TM

# Triolet Meta-Library

- Compiler "inlines" list-comprehension into `histogram`
- `histogram` has code paths for handling different loop structures
- Loop structure is known, so compiler can remove unused code paths

```
correlation xs ys =
  case IdxNest HintPar
               (arraySlice xs)
               (λx. IdxFlat HintSeq
                            (arraySlice ys)
                            (λy. f x y )          )
  of IdxNest parhint input body.
       case parhint
       of HintSeq. code for sequential nested histogram
          HintPar. parReduce input
                             (λchunk.
                                 seqHistogram 100 body chunk)
     IdxFlat parhint input body. code for flat histogram
```

# Example: Correlation Code

- Result is an outer loop specialized for this application
- Process continues for inner loop

```
correlation xs ys =
  parReduce
    (arraySlice xs)
    (λchunk. seqHistogram
            100
            (λx. IdxFlat HintSeq
                        (arraySlice ys)
                        (λy. f x y )         )
            chunk)
```

Parallel reduction; each task processes a chunk of xs
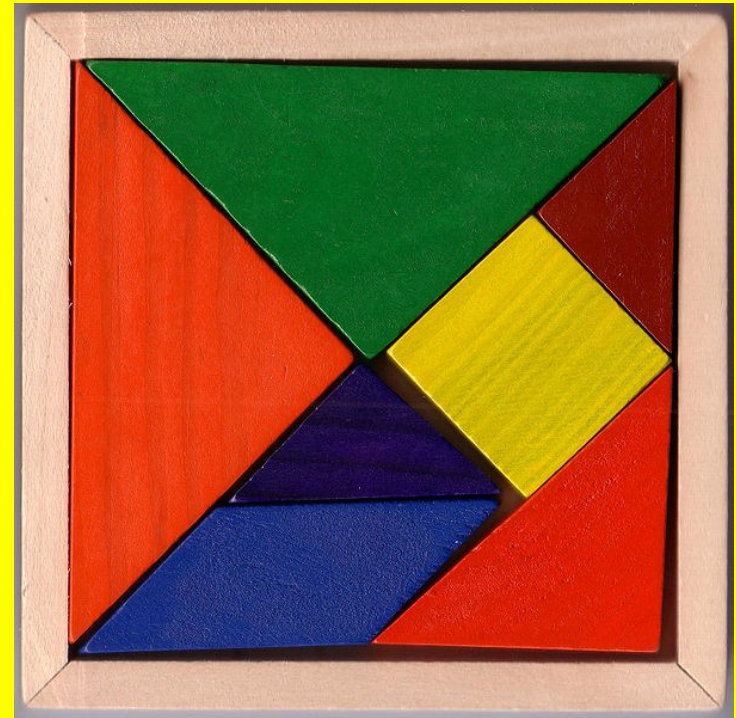
Task computes a sequential histogram

Inner loop

Body

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN
1867

# Cluster-Parallel Performance and Scalability

- Triolet delivers large speedup over sequential C

- On par with manually parallelized C for computation-bound code (left)

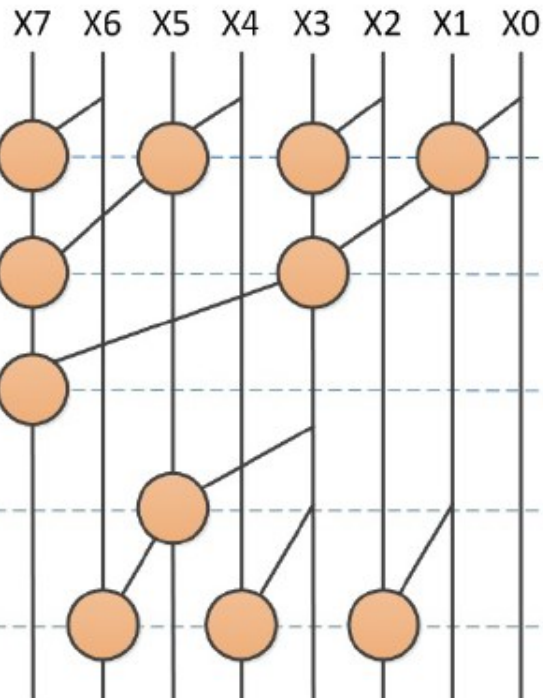- Beats similar high-level interfaces on communication-intensive code (right)

Chris Rodrigues Rodrigues, et al, PPoPP 2014

# Tangram

- A parallel algorithm framework for solving <span style="color:red">linear recurrence</span> problems
  - Scan, tridiagonal matrix solvers, bidiagonal matrix solvers, recursive filters, …
  - Many specialized algorithms in literature

- Linear Recurrence - very important for converting sequential algorithms into parallel algorithms
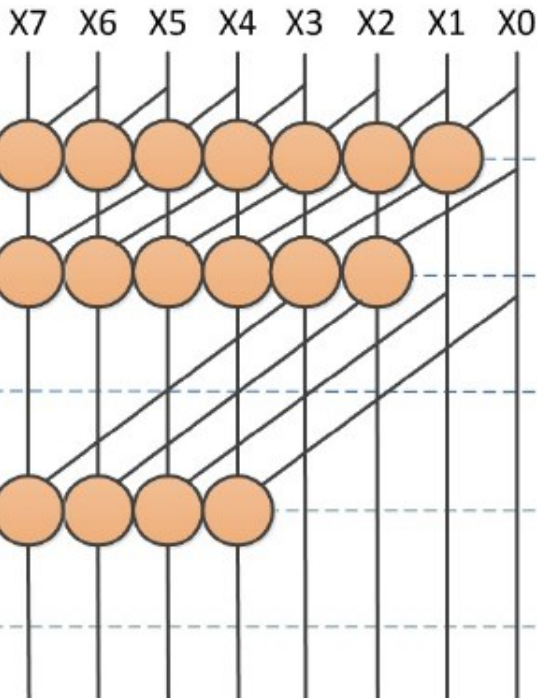
# Tangrams Linear Optimizations

- Library operations to simplify application tiling and communication
  - Auto-tuning for each target architecture
- Unified Tiling Space
  - Simple interface for register tiling, scratchpad tiling, and cache tiling
  - Automatic thread fusion as enabler
- Communication Optimization
  - Choice/hybrid of three major types of algorithms
  - Computation vs. communication tradeoff

ILLINOIS
1867
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# Linear Recurrence Algorithms and Communication



Brent-Kung Circuit

Kogge-Stone Circuit

Group Structured

Ex. Cyclic Reduction

Parallel Cyclic Reduction

Sectored Thomas

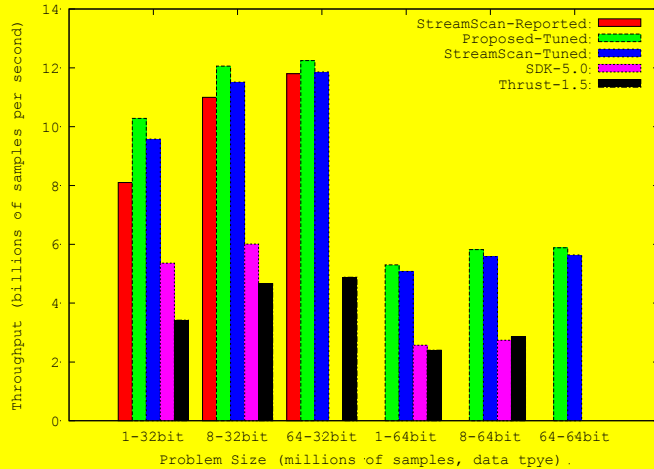# Code Programmers Need to Write: Prefix sum

- ## Sequential Code

  ```
  SEQ_Compute(…) :
  #pragma unroll
     for(…)
        UTS_REG(value,i+1) += UTS_REG(value, i);
  ```
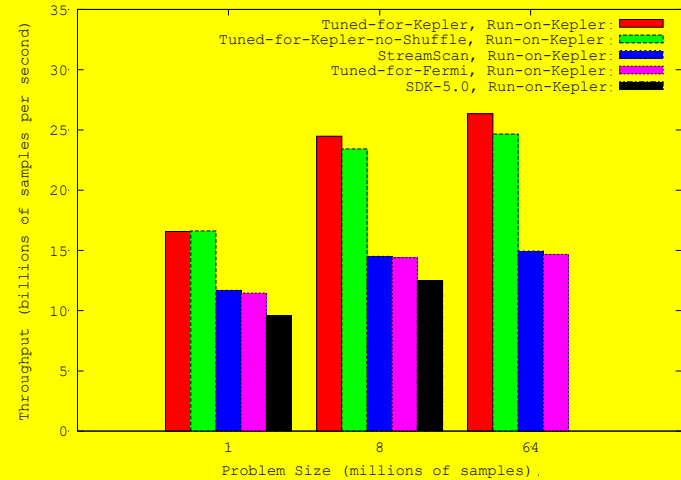
- ## Tree-Structure Code

  ```
  int p=sums[x];
  If (lane_id>=1)      sums[tx] = p = p + sums[tx-1];
  If (lane_id>=2)      sums[tx] = p = p + sums[tx-2];
  If (lane_id>=4)      sums[tx] = p = p + sums[tx-4];
  If (lane_id>=8)      sums[tx] = p = p + sums[tx-8];
  If (lane_id>=16)     sums[tx] = p = p + sums[tx-16];
  warp_sum[warp_id]=p;
  ```
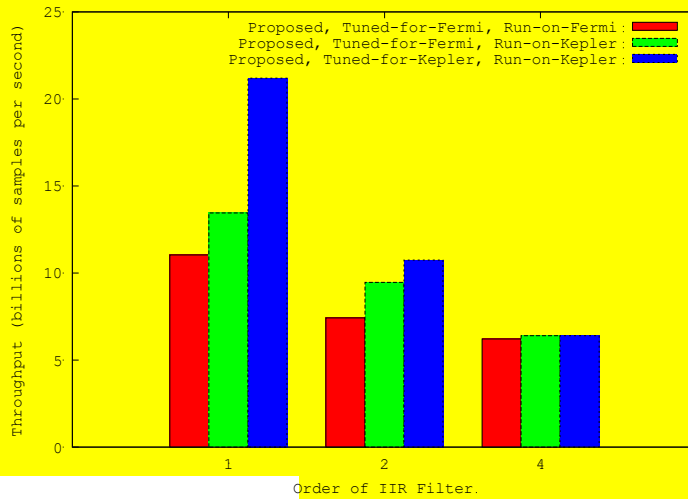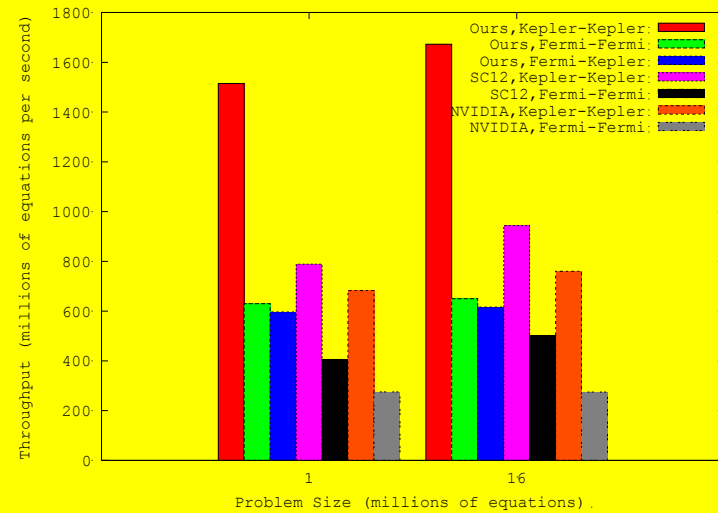
# Tangram Initial Results



Prefix scan on Fermi (C2050)

Prefix scan on Kepler(Titan)

IIR Filter on both GPUs

Tridiagonal solver on both GPUs

UNIVERSITY OF ILLINOIS

# Conclusion

- Auto-tuned generic LR algorithms in Tangram outperforms specialized scan, tridiagonal, and IIR algorithms.

- Publish and release Tangram
  - Current tridiagonal solver in CUSPARSE is from UIUC based on the Tangram work
  - Integration with Triolet

- Triolet as an open source project
  - Develop additional Triolet library functions for important application domains
  - Develop Triolet library functions for GPU targets

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN
1867
TM

**THANK YOU!**

Joint Lab 2013