# X-TUNE
# Autotuning for Exascale: Self-Tuning Software to Manage Heterogeneity

Mary Hall

April 6, 2017

THE UNIVERSITY OF UTAH

Argonne NATIONAL LABORATORY

BERKELEY LAB

X-TUNE

# Participants

| | Current | Previous |
|---|---|---|
| University of Utah | Mary Hall, Tharindu Rusira, Bob Wheeler, Derick Huth | Protonu Basu, Axel Rivera, Manu Shantharam |
| Lawrence Berkeley National Laboratory | Sam Williams, Protonu Basu | Lenny Oliker, Brian van Straalen, Phil Colella |
| Argonne National Laboratory | Prasanna Balaprakash, Paul Hovland | Thomas Nelson (Colorado), Jeff Hammond, Sri Krishna Narayanan, Stefan Wild |
| USC/ISI | | Jacqueline Chame |

THE UNIVERSITY OF UTAH    Argonne NATIONAL LABORATORY    BERKELEY LAB    X-TUNE

# Which version would you prefer to write?

/* Laplacian 7-point Variable-Coefficient Stencil */
```
for (k=0; k<N; k++)
   for (j=0; j<N; j++)
      for (i=0; i<N; i++
         temp[k][j][i] = b * h2inv * (
            beta_i[k][j][i+1] * ( phi[k][j][i+1] – phi[k][j][i] )
              -beta_i[k][j][i] * ( phi[k][j][i] – phi[k][j][i-1] )
            +beta_j[k][j+1][i] * ( phi[k][j+1][i] – phi[k][j][i] )
            -beta_j[k][j][i] * ( phi[k][j][i] – phi[k][j-1][i] )
            +beta_k[k+1][j][i] * ( phi[k+1][j][i] – phi[k][j][i] )
            -beta_k[k][j][i] * ( phi[k][j][i] – phi[k-1][j][i] ) );
```

/* Helmholz */
```
for (k=0; k<N; k++)
   for (j=0; j<N; j++)
      for (i=0; i<N; i++)
         temp[k][j][i] = a * alpha[k][j][i] * phi[k][j][i] –
                         temp[k][j][i];
```

/* Gauss-Seidel Red Black Update */
```
for (k=0; k<N; k++)
   for (j=0; j<N; j++)
      for (i=0; i<N; i++){
         if ((i+j+k+color)%2 == 0 )
         phi[k][j][i] = phi[k][j][i] – lambda[k][j][i] *
            (temp[k][j][i] – rhs[k][j][i]);}
```
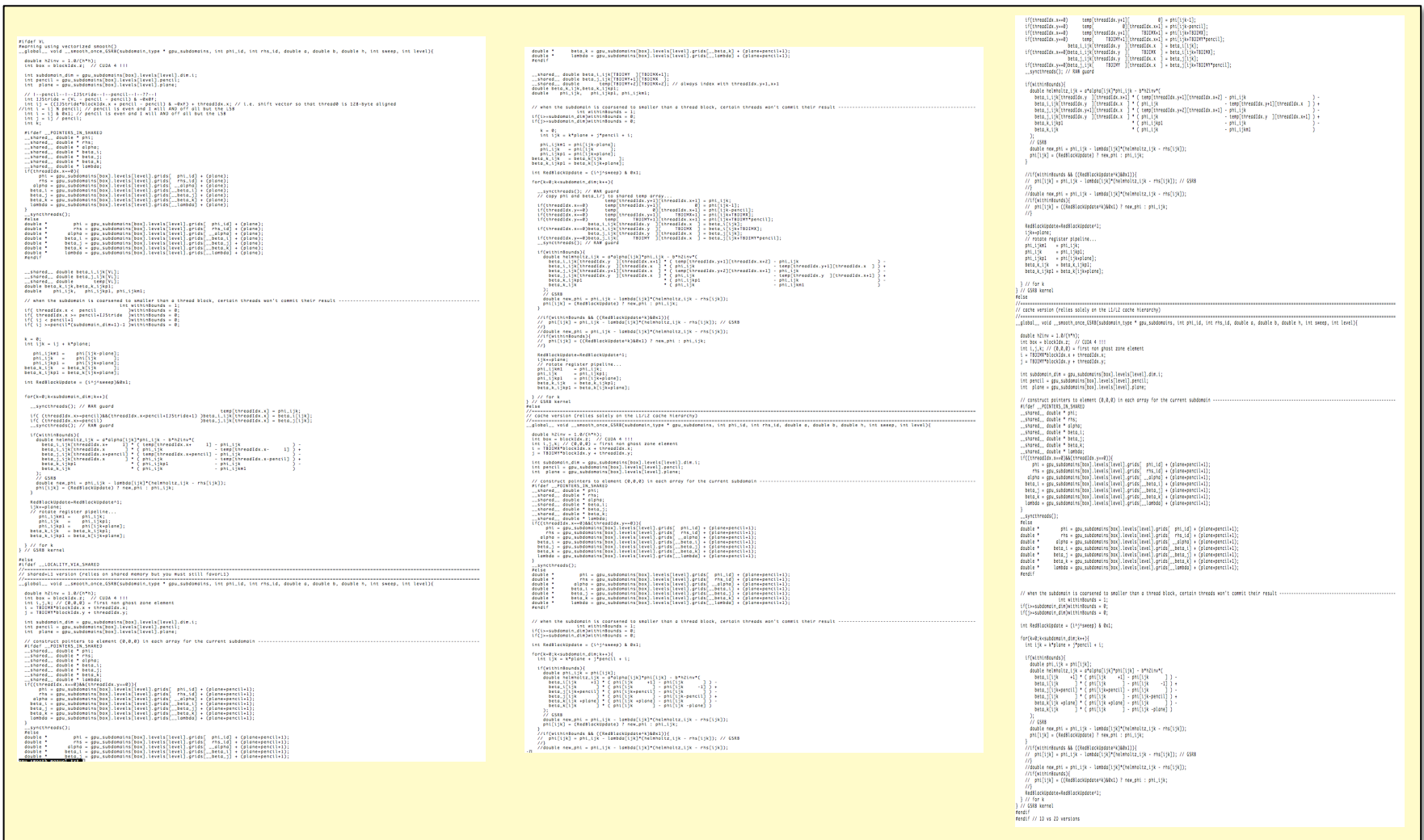
**Code A:** miniGMG baseline smooth operator approximately 13 lines of code

**Code B:** miniGMG optimized smooth operator approximately 170 lines of code

THE UNIVERSITY OF UTAH   Argonne NATIONAL LABORATORY   BERKELEY LAB   X-TUNE

# And now GPU code?



**Code C:** miniGMG optimized smooth operator for GPU, 308 lines of code for just kernel

THE UNIVERSITY OF UTAH · Argonne NATIONAL LABORATORY · BERKELEY LAB · X-TUNE

# Which version would you prefer to write?

/* local_grad_3 computation from nek5000 */

w[nelt i j k] += Dt[l k] U[nelt n m l] D[j m] D[i n]

## Code A:

1 line mathematical representation
Input to OCTOPI

/* local_grad3 from nek5000, generated CUDA code */



## Code B:

Generated CUDA+harness, 122 lines of code

# Exascale Challenges: Code B is not Unusual

- Performance portability?
  - Particularly across fundamentally different CPU and GPU architectures

- Programmer productivity?
  - High performance implementations will require low-level specification that exposes architecture

- Software maintainability and portability?
  - May require multiple implementations of application

## Current solutions

- Follow MPI and OpenMP standards
  - Same code unlikely to perform well across CPU and GPU
  - Vendor C and Fortran compilers not optimized for HPC workloads

- Some domain-specific framework strategies
  - Libraries, C++ template expansion, standalone DSL
  - Not *composable* with other optimizations

THE UNIVERSITY OF UTAH

Argonne NATIONAL LABORATORY

BERKELEY LAB

X-TUNE

# Programming Systems Must Hide Complexity!

- Define a common abstraction(s) that programmers can target

- Use compilation tools to map to optimal implementation
  - hide **programming model choices** from users
  - hide **architectural complexity** from users
  - hide **tuning** from users

# X-TUNE Goals

- Application programmer expresses key computation at a high level *(Code A)*
  - Sequential C or domain-specific specification
- Code transformations are applied
  - Existing and domain-specific transformations
  - Generates a collection of optimized implementations
  - Includes thread-level code generation (OpenMP and CUDA)
- Autotuning
  - Searches the space of implementations to find the best match to execution context
  - Selects optimized implementation *(Code B)*
- Automation mitigates correctness, productivity, portability, maintainability concerns

*X-TUNE automates the process of converting Code A to Code B.*
*See today's demonstrations!*

# X-TUNE Approach

- For each *motif*, start with manually-tuned code or work with developer of new code
  - What transformations are needed to target specific architectures?
  - What performance questions can be addressed by autotuning?

- Attempt to automate
  - Exploit existing compiler transformations
  - Develop new domain-specific transformations and required analysis and code generation support
  - Develop decision algorithms

- Collect application code from collaborators, Co-Design Centers and other DOE application teams
  - Generalize from experiments with manually-tuned code

# X-TUNE Key Ideas



**CUDA-CHiLL**

- **Composable** transformation and code generation
  - Leverage rich set of existing transformations and code generation capability in **polyhedral framework**

- Extensible to domain-specific transformations and decision algorithms
  - **Compose** with existing transformations

- Optimization strategies and parameters exposed to autotuning via transformation recipes

- Search space navigation
  - Search framework can be standalone tool (e.g., Orio, OpenTuner, Active Harmony, Nitro)

THE UNIVERSITY OF UTAH · Argonne NATIONAL LABORATORY · BERKELEY LAB · **X-TUNE**

# Example Motifs Supported by X-TUNE

| Motif | Input | Existing Transformations | Domain-specific transformations | Autotuning | Search | |
|---|---|---|---|---|---|---|
| Geometric Multigrid | Sequential C computation (w/ MPI and OpenMP harness) | Communication-avoiding: fusion, tile, wavefront (skew&permute), OpenMP, CUDA | Ghost zones, Partial sums | Ghost zone depth, threading, strategy at each level of V-cycle | Simple, full space | X-TUNE DEMO |
| Tensor Contraction | Mathematical Formula | Tile, permute, scalar replacement, unroll, CUDA | Rewriting, Decision algorithm | Loop order, CUDA threading | SURF | X-TUNE DEMO |
| Sparse Matrix Computation | Sequential C with CSR matrix | Tile, permute, skew, unroll, reduction, scalar expansion, OpenMP, CUDA | Generate inspectors, coalesce, make-dense, compact, split, level sets | Threading, matrix repr. | Simple, full space | SUPER/NSF |

THE UNIVERSITY OF UTAH  ·  Argonne NATIONAL LABORATORY  ·  BERKELEY LAB  ·  X-TUNE

# Geometric Multigrid

- Multigrid solves elliptic PDEs in O(N) computational complexity by using a hierarchical approach.

smooth $Lu^h = f^h$

$r^h = f^h - Lu^h$ (residual)
$f^{2h} = restrict(r^h)$

smooth $Lu^{2h} = f^{2h}$

$r^{2h} = f^{2h} - Lu^{2h}$
$f^{4h} = restrict(r^{2h})$

smooth $Lu^{4h} = f^{4h}$

$r^{4h} = f^{4h} - Lu^{4h}$
$f^{8h} = restrict(r^{4h})$

multiple smooth's on $Lu^{8h} = f^{8h}$
(or Iterative Solver like BiCGStab)

smooth $Lu^h = f^h$

$u^h$ += interpolate($u^{2h}$)

smooth $Lu^{2h} = f^{2h}$

$u^{2h}$ += interpolate($u^{4h}$)

smooth $Lu^{4h} = f^{4h}$

$u^{4h}$ += interpolate($u^{8h}$)

*progress within V-cycle*

- ❖ Geometric Multigrid (**GMG**) is specialization in which the operator (A) is simply a stencil on a structured grid (i.e. *matrix-free*)
- ❖ **Stresses performance at different degrees of parallelism, locality, working set sizes, etc…**
- ❖ **Optimization strategy varies across different levels of V-cycle, even on one architecture!**

THE UNIVERSITY OF UTAH
Argonne NATIONAL LABORATORY
BERKELEY LAB
X-TUNE

12

# miniGMG Benchmark

- **miniGMG proxies the MG solves in BoxLib/Chombo codes**
  - **Predecessor to ExACT Co-Design Center's HPGMG**
- Distributed memory (MPI) implementation
- **operator** is configurable
  - 7pt variable coefficient **proxies LMC**
  - 7pt constant coefficient is simpler
  - 125pt/27pt/13pt high-order stencils
- **smoother** in the v-cycle is configurable
  - Gauss Seidel, Red-Black (GSRB) = **proxies**
  - Jacobi (mathematically weaker)
- **bottom solver** is configurable
  - multiple GSRB's
  - Krylov solver like **BiCGStab**, CG, CA-BiCGStab, CA-CG, etc...

one subdomain of $64^3$ elements

Collection of subdomains owned by an MPI process

THE UNIVERSITY OF UTAH

Argonne NATIONAL LABORATORY

BERKELEY LAB

X-TUNE

# Optimized Code A can beat Code B!

- miniGMG optimized w/CHiLL
  - fused operations
  - created a communication-avoiding wavefront
  - **auto-parallelized (OpenMP)**
- **autotuning** finds the best implementation for each box size
  - wavefront depth (degree of comm. avoiding)
  - Turn on/off optimizations (fusion, wavefront)
  - nested OpenMP configuration
  - inter-thread synchronization (barrier vs. P2P)
- For fine grids (large arrays) CHiLL attains a **4.3x speedup** over baseline

## GSRB Smooth (Edison)



Legend:
- CHiLL generated
- Manually Tuned
- Baseline

Y-axis: Speedup over Baseline Smoother (0.0x – 5.0x)

X-axis: Box Size ( == Level in the V-Cycle) — 64^3, 32^3, 16^3, 8^3, 4^3

THE UNIVERSITY OF UTAH · Argonne NATIONAL LABORATORY · BERKELEY LAB · X-TUNE

# Flexibility

- Fuse the residual and restriction operations into the wavefront as well
  - read $u^h$, $R^h$, and coefficients once
  - perform 4 smooths (**no additional data movement**)
  - write smoothed $u^h$ and new $R^{2h}$

- Apply these transformations to a different smoother and autotune it
  - up to **3x improvement in MGSolve**



X-TUNE

# Retargetable and Performance Portable: Optimized Code A can beat Code C!

- CHiLL can obviate the need for architecture-specific programming models like CUDA

  - CUDA-CHiLL took the sequential GSRB implementation (.c) and **generated CUDA** that runs on NVIDIA GPUs

  - CUDA-CHiLL tunes for the current target machine whereas static implementations hand-optimize for yesterday's GPUs

  - CUDA-CHiLL autotuned over the thread block sizes and is ultimately **2% faster** than the hand-optimized minigmg-cuda (**Code C**)

GSRB Smooth on 64^3 boxes

Time (seconds) vs 2D Thread Blocks <TX,TY>

Legend: CUDA-CHiLL, Handtuned, Handtuned-VL

5.224148    4.861889    4.774941

THE UNIVERSITY OF UTAH    Argonne NATIONAL LABORATORY    BERKELEY LAB    X-TUNE

# Extensible to Domain-Specific Optimizations

- Applied mathematicians are exploring how changing the stencils may be better suited for future architecture trends
- Consider the following variations (stencils) on the discretization of the Laplacian
  - low-level implementations (optimized OMP4) may provide high performance
  - but are one-off solutions as requisite optimizations/tuning change from one stencil to the next

THE UNIVERSITY OF UTAH

Argonne NATIONAL LABORATORY

BERKELEY LAB

X-TUNE

# Extensible to Domain-Specific Optimizations

■ CHiLL optimized/tuned each of these stencils

- Introduced **partial sums** optimization to avoid redundant computation for compute-bound high-order stencils

- selected unique optimizations for each stencil and at each level of the MG V-Cycle

- Without a communication-avoiding wavefront, CHiLL delivered performance **near the Roofline bound**.

- Using a wavefront, CHiLL can nearly **double** the nominal Roofline performance for the 7- and 27-point operators.



Smoother Performance (Fine Grid)

Legend:
- All Optimizations
- +Fusion & Wavefront
- +Fusion & Partial Sums
- +Fusion
- Baseline
- Roofline Memory Bound

Edison: 7pt, 27pt, 13pt, 125pt
Hopper: 7pt, 27pt, 13pt, 125pt

# Example Transformation Recipes

- These can be manually-written (miniGMG) or automatically generated (tensor contraction)

```
/* jacobi_box_4_64.py, 27-pt stencil, 64³ box size */
from chill import *

#select which computation to optimize
source('jacobi_box_4_64.c')
procedure('smooth_box_4_64')
loop(0)
original() # fuse wherever possible

#create a parallel wavefront
skew([0,1,2,3,4,5],2,[2,1])
permute([2,1,3,4])

#partial sum for high order stencils and fuse result
distribute([0,1,2,3,4,5],2)
stencil_temp(0)
stencil_temp(5)
fuse([2,3,4,5,6,7,8,9],1)
fuse([2,3,4,5,6,7,8,9],2)
fuse([2,3,4,5,6,7,8,9],3)
fuse([2,3,4,5,6,7,8,9],4)
```

```
/* gsrb.lua, variable coefficient GSRB, 64³ box size */
init("gsrb_mod.cu", "gsrb",0,0)
dofile("cudaize.lua") # custom commands in lua

# set up parallel decomposition, adjust via autotuning
TI=32
TJ=4
TK=64
TZ=64

tile_by_index(0, {"box","k","j", "i"},{TZ,TK, TJ, TI},
{l1_control="bb", l2_control="kk", l3_control="jj",
l4_control="ii"},{"bb","box","kk","k","jj","j","ii","i"})

cudaize(0, "kernel_GPU",
{_temp=N*N*N*N,_beta_i=N*N*N*N,
_phi=N*N*N*N},{block={"ii","jj","box"},
thread={"i","j"}},{})
```

# Tensor Contraction: Spectral Element Method from nek5000/nekbone (CESAR)

$$C = A \otimes B \underline{u}$$

- A and B are square matrices
- $\underline{u}$ is a component vector
- In 2-d, C can be computed:

$$c_{i,j} = \sum_l \sum_k a_{j,l} b_{i,k} u_{k,l}$$

**Order O(n⁴)**

## Optimize by rewriting to the following:

$$C = (A \otimes I)(I \otimes B)\underline{u}$$

Partial Results: $\underline{w} = (I \otimes B)\underline{u}$ $\longrightarrow$ $w_{i,j} = \sum_l u_{i,l} b_{l,j}^T$ **Order O(n³)**

Final Results: $C = (A \otimes I)\underline{w}$ $\longrightarrow$ $c_{i,j} = \sum_k a_{i,k} w_{k,j}$

THE UNIVERSITY OF UTAH   Argonne NATIONAL LABORATORY   BERKELEY LAB   X-TUNE

# Tensor Contraction: Challenging Mapping, Particularly for GPU

- What is the optimal contraction order?
- What is the optimal loop order? (N! different implementations)
- GPU challenges: small dimensions, memory hierarchy effects
- Search space is discontinuous, noisy, and expensive to evaluate
- X-TUNE Approach:
  - Fully automate from mathematical description to GPU code generation (**Code A** to **Code B**)
  - Automate and reduce search time across intractable brute force search space

THE UNIVERSITY OF UTAH   Argonne NATIONAL LABORATORY   BERKELEY LAB   **X-TUNE**

# Baracuda Framework



**Key points:**

CHiLL as an Orio module provides search
space and code generation
SURF manages exploration of search space

# SURF: Model-Based Search

- **Search Using Random Forests**
  - state-of-the-art statistical machine learning algorithm
  - handles binary permutation parameters
  - handles nonlinear parameter interactions
- Approach
  - start with promising small set of parameter configurations
  - evaluate performance
  - fit surrogate model (ML)
  - predict new set of high-performing configurations
  - iterate



*Example Surrogate Model Fitted to Sampled Performance
(iterative refinement improves the statistical model)*

# Optimizing NWChem

- Extracted representative on-node tensor contractions from NWChem/TCE
  - many small contractions
  - atypical of OpenACC use model
- Baracuda generates optimized CUDA for NVIDIA's Fermi or Kepler GPUs
- Manually modified CUDA to OpenACC
  - Naïve replaces CUDA with OpenACC, but uses same loop order and parallelization
  - OpenACC = naïve + manual explicit control over hierarchical parallelism



X-TUNE

# Optimizing Nekbone

- Nekbone (CESAR CoDesign Center Proxy App) with optimized local_grad_3 and local_grad_3t
  - Many, small (e.g. 12x12x12) contractions
  - Nominally implemented as many BLAS3 calls
- Baracuda generates optimized CUDA for NVIDIA's Fermi or Kepler GPUs
- Compare to single Haswell core.

| Speedup over 1-core Haswell | Naïve OpenACC | Optimized OpenACC | Baracuda (CUDA) |
|---|---|---|---|
| K20 | 2.86 | 12.39 | 36.47 |
| C2050 | 1.18 | 19.21 | 34.65 |

THE UNIVERSITY OF UTAH    Argonne NATIONAL LABORATORY    BERKELEY LAB    X-TUNE

# Summary of X-TUNE Accomplishments

- Demonstrated for two motifs, Geometric Multigrid and Tensor Contraction
    - Automated architecture-specific optimization from high-level specification
    - Performance rivaling manually-tuned code and sometimes better
    - Approach can achieve performance portability, productivity and maintainability

- Implementation status
    - GMG optimizations integrated into CHiLL, tensor contraction uses existing CHiLL with additional frontend and Orio
    - CHiLL publicly available on github
    - Installed on Edison (user space)
    - Demonstrations this evening

**X-TUNE**

# Publications

- Papers
    - P. Basu, M. Hall, M. Khan, S.Maindola, S.Muralidharan, S.Ramalingam, A.Rivera, M.Shantharam, A.Venkat. Towards Making Autotuning Mainstream. International Journal of High Performance Computing Applications, 27(4), November 2013.
    - P. Basu, S. Williams, A. Venkat, B. Van Straalen, M. Hall, and L. Oliker. Compiler generation and autotuning of communication-avoiding operators for geometric multigrid. In High Performance Computing Conference (HIPC), 2013.
    - P. Basu, S. Williams, A. Venkat, B. Van Straalen, M. Hall, and L. Oliker. Compiler generation and autotuning of communication-avoiding operators for geometric multigrid. In Workshop on Optimizing Stencil Computations (WOSC), 2013.
    - Protonu Basu, Samuel Williams, Brian Van Straalen, Mary Hall, Leonid Oliker, Phillip Colella, "Compiler-Directed Transformation for Higher-Order Stencils", International Parallel and Distributed Processing Symposium (IPDPS), May 2015.
    - Thomas Nelson, Axel Rivera, Prasanna Balaprakash, Mary Hall, Paul D. Hovland, Elizabeth Jessup, Boyana Norris, "Generating Efficient Tensor Contractions for GPUs", International Conference on Parallel Processing (ICPP), September 2015.

- Thesis and Dissertations
    - Axel Rivera. Using Autotuning for Accelerating Tensor-Contraction on GPUs, Masters thesis, University of Utah, December 2014.
    - Protonu Basu, "Compiler Optimizations and Autotuning for Stencils and Geometric Multigrid", PhD Dissertation, University of Utah, December 2015.
    - Thomas Nelson, "DSLs and Search for Linear Algebra Performance Optimization," PhD Dissertation, University of Colorado, December 2015.

THE UNIVERSITY OF UTAH    Argonne NATIONAL LABORATORY    BERKELEY LAB    X-TUNE