

The Pluto Compiler and its Use for Computational Sciences

Uday Bondhugula

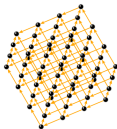
uday@csa.iisc.ernet.in

Dept of CSA

Indian Institute of Science

Bangalore 560012 India

**with Aravind Acharya, Vinay Vasista, Kumudha KN, Irshad
Pananilath, Ravi Teja Mullapudi**



WHAT IS PLUTO?

- **A source-to-source optimizer and parallelizer**

WHAT IS PLUTO?

- **A source-to-source optimizer and parallelizer**
- **Uses many other polyhedral libraries and tools like ISL, Polylib, Cloog, Pet, Clan, Candi**

HOW CAN PLUTO BE USED?

- **Push button:** fully automatically for optimization (tiling and other transformations), parallelization
- **Almost automatic:** With an understanding of what Pluto does, use it to obtain desired result
- **DSLs:** In domain-specific compilers/optimizers or library generators

HOW CAN PLUTO BE USED?

- **Push button:** fully automatically for optimization (tiling and other transformations), parallelization
- **Almost automatic:** Having understood what Pluto does, use it to obtain desired result
- **DSLs:** In domain-specific compilers/optimizers or library generators

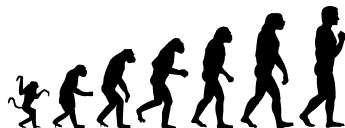
BIG PICTURE: ROLE OF COMPILERS

- Improve existing **general-purpose** compilers (for C, C++, Python, ...)
- LLVM/Polly, GCC/Graphite, PPCG, Pluto, other compilers/tools
- Build new **domain-specific languages and compilers**
- Scientists say **WHAT** they execute and not **HOW** they execute

BIG PICTURE: ROLE OF COMPILERS

The Evolutionary Approach

- Improve existing **general-purpose** compilers (for C, C++, Python, ...)
- LLVM/Polly, GCC/Graphite, PPCG, Pluto, other compilers/tools



The Revolutionary Approach

- Build new **domain-specific languages and compilers**
- Scientists say **WHAT** they execute and not **HOW** they execute



Important to pursue both

- 1 Pluto/Pluto+
 - Affine Transformations
 - Tiling

- 2 Case Studies
 - Solving Partial Differential Equations
 - Lattice Boltzmann Method
 - Image Processing Pipelines

- 3 Conclusions

- 1 Pluto/Pluto+
 - Affine Transformations
 - Tiling
- 2 Case Studies
 - Solving Partial Differential Equations
 - Lattice Boltzmann Method
 - Image Processing Pipelines
- 3 Conclusions

AFFINE TRANSFORMATIONS

- Examples of affine functions of i, j : $i + j, i - j, i + 1, 2i + 5$
- Not affine: $ij, i^2, i^2 + j^2, a[j]$

AFFINE TRANSFORMATIONS

- Examples of affine functions of i, j : $i + j, i - j, i + 1, 2i + 5$
- Not affine: $ij, i^2, i^2 + j^2, a[j]$

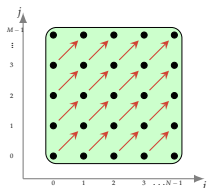


Figure: Iteration space

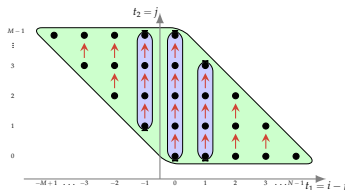


Figure: Transformed space

```
for (i = 0; i < N; i++){
  for (j = 0; j < M; j++){
    A[i+1][j+1] = f(A[i][j])
  }
}
```

```
#pragma omp parallel for private(t2)
for (t1=-M+1; t1<=N-1; t1++) {
  for (t2=max(0, -t1); t2<=min(M-1, N-1-t1); t2++){
    A[t1+t2+1][t2+1] = f(A[t1+t2][t2]);
  }
}
```

- Transformation: $(i, j) \rightarrow (i - j, j)$

AFFINE TRANSFORMATIONS

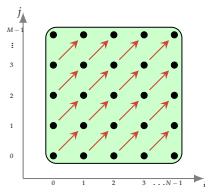


Figure: Iteration space

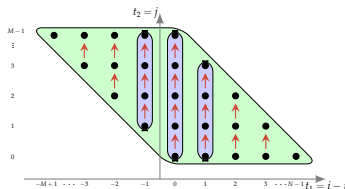


Figure: Transformed space

- Affine transformations are attractive because:
 - Preserve **collinearity** of points and **ratio of distances** between points
 - Code generation with affine transformations has thus been studied well (CLooG, ISL, OMEGA+)

AFFINE TRANSFORMATIONS

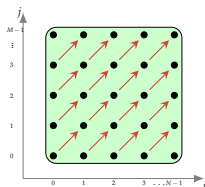


Figure: Iteration space

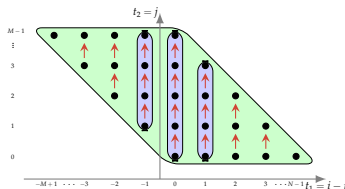


Figure: Transformed space

- Affine transformations are attractive because:
 - Preserve **collinearity** of points and **ratio of distances** between points
 - Code generation with affine transformations has thus been studied well (CLooG, ISL, OMEGA+)
 - Model a very rich class of loop re-orderings
 - Useful for several domains like dense linear algebra, stencils, image processing pipelines, Lattice Boltzmann Method

AFFINE TRANSFORMATIONS

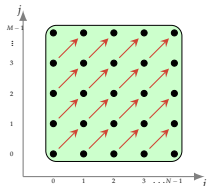


Figure: Iteration space

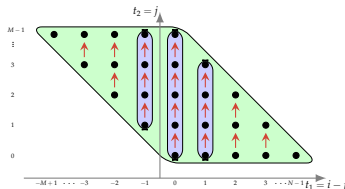


Figure: Transformed space

- Affine transformations can improve **parallelism** and **locality** (Feautrier 1992, Lengauer, Lim and Lam 1997, Griehl 2004, Pluto 2008)

- Designed around 2008 [Bondhugula et al. CC 2008, PLDI 2008]
- **Finds good transformations to improves locality and parallelism**
- Extended in 2014-2015 (transformation coefficients need not be non-negative)

FINDING VALID AND GOOD AFFINE TRANSFORMATIONS

(i, j)

(j, i)

$(i+j, j)$

$(i-j, j)$

$(i, i+j)$

$(i+j, i-j)$

...

FINDING VALID AND GOOD AFFINE TRANSFORMATIONS

(i, j)
 (j, i)
 $(i+j, j)$
 $(i-j, j)$
 $(i, i+j)$
 $(i+j, i-j)$
...

- One-to-one functions
- Validity: dependences should not be violated
- Coefficients: for $i - j$, the coefficients are 1,-1

- Optimization Problem: **Minimize dependence distance**

- Optimization Problem: **Minimize dependence distance**
- Constraints:
 - Tiling validity constraints
 - Dependence distance bounding constraints
 - Linear Independence constraints

- Optimization Problem: **Minimize dependence distance**
- Constraints:
 - Tiling validity constraints
 - Dependence distance bounding constraints
 - Linear Independence constraints

TILING (BLOCKING)

- Partition and execute iteration space in blocks

```
for (i=1; i<T; i++)  
  for (j=1; j<N-1; j++)  
    S(i,j)
```

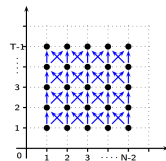


Figure: Iteration space

TILING (BLOCKING)

- Partition and execute iteration space in blocks
- Benefits - *cache locality & parallelism*

```
for (i=1; i<T; i++)  
  for (j=1; j<N-1; j++)  
    S(i,j)
```

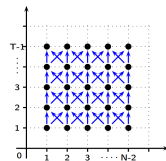


Figure: Iteration space

TILING (BLOCKING)

- Partition and execute iteration space in blocks
- Benefits - *cache locality & parallelism*
- Validity of tiling

```
for (i=1; i<T; i++)  
  for (j=1; j<N-1; j++)  
    S(i,j)
```

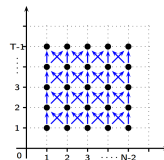


Figure: Iteration space

TILING (BLOCKING)

- Partition and execute iteration space in blocks
- Benefits - *cache locality & parallelism*
- Validity of tiling
 - No cycle between tiles

```
for (i=1; i<T; i++)  
  for (j=1; j<N-1; j++)  
    S(i,j)
```

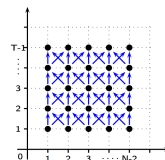


Figure: Iteration space

TILING (BLOCKING)

- Partition and execute iteration space in blocks
- Benefits - *cache locality & parallelism*
- Validity of tiling
 - No cycle between tiles
 - Sufficient condition: All dependence components should be non-negative

```
for (i=1; i<T; i++)  
  for (j=1; j<N-1; j++)  
    S(i,j)
```

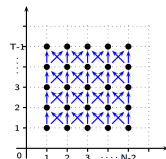


Figure: Iteration space

TILING (BLOCKING)

- Partition and execute iteration space in blocks
- Benefits - *cache locality & parallelism*
- Validity of tiling
 - No cycle between tiles
 - Sufficient condition: All dependence components should be non-negative
- Time tiling

```
for (i=1; i<T; i++)  
  for (j=1; j<N-1; j++)  
    S(i,j)
```

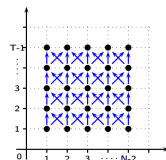


Figure: Iteration space

TILING (BLOCKING)

- Partition and execute iteration space in blocks
- Benefits - *cache locality & parallelism*
- Validity of tiling
 - No cycle between tiles
 - Sufficient condition: All dependence components should be non-negative
- Time tiling

```
for (i=1; i<T; i++)  
  for (j=1; j<N-1; j++)  
    S(i,j)
```

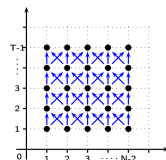


Figure: Iteration space

TILING (BLOCKING)

- Partition and execute iteration space in blocks
- Benefits - *cache locality & parallelism*
- Validity of tiling
 - No cycle between tiles
 - Sufficient condition: All dependence components should be non-negative
- Time tiling

```
for (i=1; i<T; i++)  
  for (j=1; j<N-1; j++)  
    S(i,j)
```

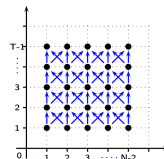


Figure: Iteration space

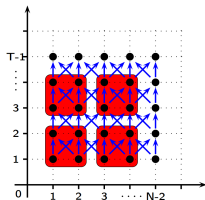


Figure: Invalid tiling

TILING (BLOCKING)

- Partition and execute iteration space in blocks
- Benefits - *cache locality & parallelism*
- Validity of tiling
 - No cycle between tiles
 - Sufficient condition: All dependence components should be non-negative
- Time tiling

```
for (i=1; i<T; i++)  
  for (j=1; j<N-1; j++)  
    S(i,j)
```

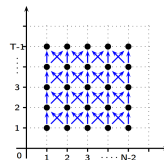


Figure: Iteration space

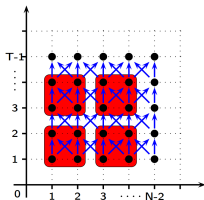


Figure: Invalid tiling

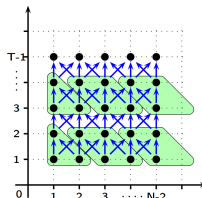


Figure: Valid tiling

CONCURRENT START AND DIAMOND TILING¹

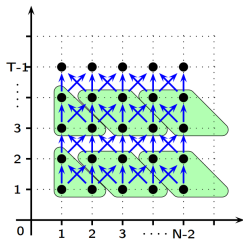


Figure: Parallelogram tiling

CONCURRENT START AND DIAMOND TILING¹

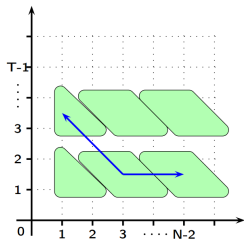


Figure: Pipelined start

CONCURRENT START AND DIAMOND TILING¹

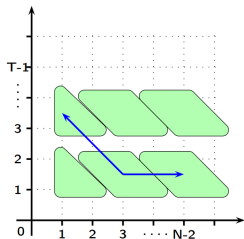


Figure: Pipelined start

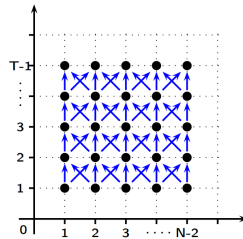


Figure: Original iteration space

CONCURRENT START AND DIAMOND TILING¹

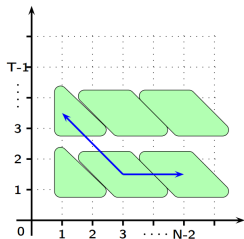


Figure: Pipelined start

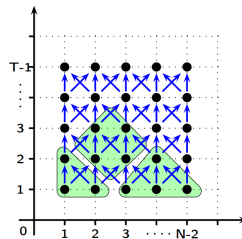


Figure: Group as diamonds

CONCURRENT START AND DIAMOND TILING¹

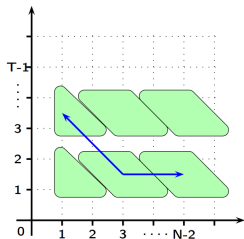


Figure: Pipelined start

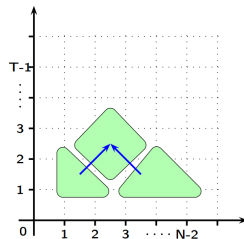


Figure: Concurrent start possible

CONCURRENT START AND DIAMOND TILING¹

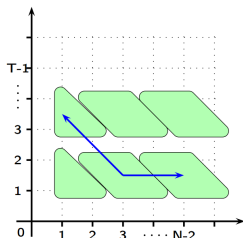


Figure: Pipelined start

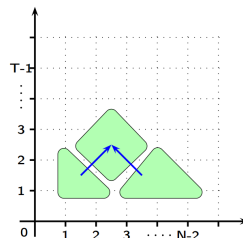


Figure: Concurrent start possible

- Diamond tiling

- Face allowing concurrent should be strictly within the cone of the tiling hyperplanes
- Eg: $(1,0)$ is in the cone of $(1,1)$ and $(1,-1)$

CLASSICAL TIME SKEWING VS DIAMOND TILING

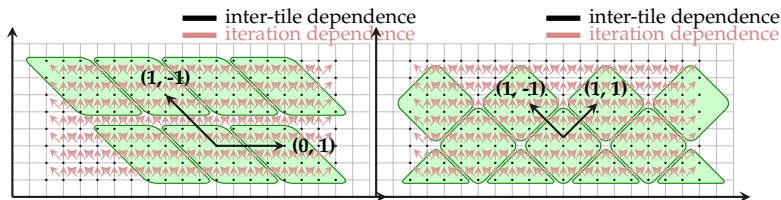


Figure: Two ways of tiling heat-1d: parallelogram & diamond

- Classical time skewing: $(t, i) \rightarrow (t, t + i)$
- Diamond tiling: $(t, i) \rightarrow (t + i, t - i)$

A SEQUENCE OF TRANSFORMATIONS FOR 2-D JACOBI RELAXATIONS

```
for (t = 0; t < T; t++)  
  for (i = 1; i < N+1; i++)  
    for (j = 1; j < N+1; j++)  
      A[(t+1)%2][i][j] = f((A[t%2][i+1][j], A[t%2][i][j], A[t%2][i-1][j],  
                           A[t%2][i][j+1], A[t%2][i][j-1], A[t%2][i][j]));
```

- 1 Enabling transformation for diamond tiling

$$T((t, i, j)) = (t + i, t - i, t + j).$$

A SEQUENCE OF TRANSFORMATIONS FOR 2-D JACOBI RELAXATIONS

```
for (t = 0; t < T; t++)  
  for (i = 1; i < N+1; i++)  
    for (j = 1; j < N+1; j++)  
      A[(t+1)%2][i][j] = f((A[t%2][i+1][j], A[t%2][i][j], A[t%2][i-1][j],  
                           A[t%2][i][j+1], A[t%2][i][j-1], A[t%2][i][j]));
```

- 1 Enabling transformation for diamond tiling

$$T((t, i, j)) = (t + i, t - i, t + j).$$

- 2 Perform the actual tiling (in the transformed space)

$$T'((t, i, j)) = \left(\frac{t + i}{64}, \frac{t - i}{64}, \frac{t + j}{64}, t + i, t - i, t + j \right)$$

A SEQUENCE OF TRANSFORMATIONS FOR 2-D JACOBI RELAXATIONS

```
for (t = 0; t < T; t++)  
  for (i = 1; i < N+1; i++)  
    for (j = 1; j < N+1; j++)  
      A[(t+1)%2][i][j] = f((A[t%2][i+1][j], A[t%2][i][j], A[t%2][i-1][j],  
        A[t%2][i][j+1], A[t%2][i][j-1], A[t%2][i][j]));
```

- 1 Enabling transformation for diamond tiling

$$T((t, i, j)) = (t + i, t - i, t + j).$$

- 2 Perform the actual tiling (in the transformed space)

$$T'((t, i, j)) = \left(\frac{t + i}{64}, \frac{t - i}{64}, \frac{t + j}{64}, t + i, t - i, t + j \right)$$

- 3 Create a wavefront of tiles

$$T''((t, i, j)) = \left(\frac{t + i}{64} + \frac{t - i}{64}, \frac{t - i}{64}, \frac{t + j}{64}, t, t + i, t + j \right)$$

SOME EXAMPLES OF HOW PLUTO CAN BE USED

- Optimize Jacobi and other relaxations via time tiling

SOME EXAMPLES OF HOW PLUTO CAN BE USED

- Optimize Jacobi and other relaxations via time tiling
- Optimize pre-smoothing steps at various levels of Geometric Multigrid method

SOME EXAMPLES OF HOW PLUTO CAN BE USED

- Optimize Jacobi and other relaxations via time tiling
- Optimize pre-smoothing steps at various levels of Geometric Multigrid method
- Optimize Lattice Boltzmann Method computations

Web: <http://pluto-compiler.sf.net>

- Use git version
- Use 'pet' branch of git version
- Preferable: use Intel's C/C++ compiler (14.0 or higher) to compile generated code

- 1 Pluto/Pluto+
 - Affine Transformations
 - Tiling
- 2 Case Studies
 - Solving Partial Differential Equations
 - Lattice Boltzmann Method
 - Image Processing Pipelines
- 3 Conclusions

POISSON'S EQUATION

Poisson's equation:

$$\nabla^2 u = f.$$

$$\frac{1}{h^2} \begin{bmatrix} & -1 & \\ -1 & 4 & -1 \\ & -1 & \end{bmatrix} u_h = f_h$$

- We are solving $y = Ax$
- What about A^{-1} ?

- Use a hierarchical structure – a multi-scale representation of the grid
- Perform pre-smoothing at a finer level
- Restrict the error to a coarser grid
- Solve for the error at a coarser level (recursion)
- Interpolate the error to the finer level

- Run multiple iterations of the above

Pluto can be used to optimize the pre-smoothing or post-smoothing steps readily

HIERARCHICAL MESH STRUCTURE

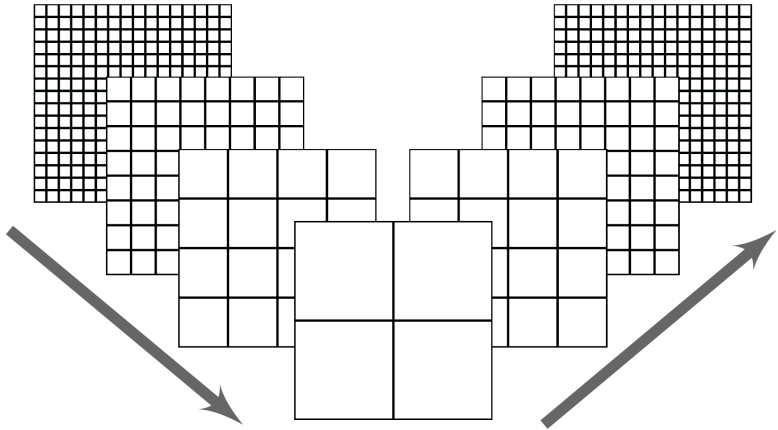
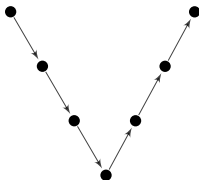


Figure: Hierarchical mesh structure for Multigrid levels

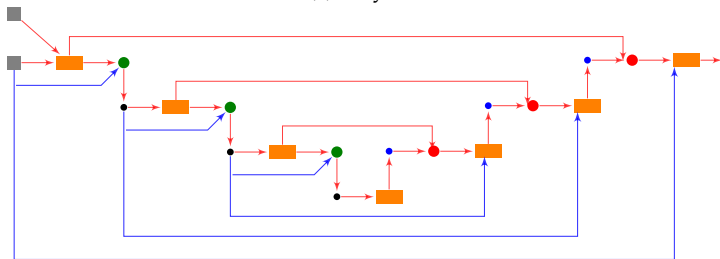
MULTIGRID V-CYCLE: ALGORITHM

```
Input :  $v^h, f^h$ 
1 Relax  $v^h$  for  $n_1$  iterations:  $v^h \leftarrow (1 - \omega D^{-1} A^h) v^h + \omega D^{-1} f^h$ 
   // pre-smoothing
2 if coarsest level then
3   | Relax  $v^h$  for  $n_2$  iterations                                // coarse smoothing
4    $r^h \leftarrow f^h - A^h v^h$                                      // residual
5    $r^{2h} \leftarrow I_h^{2h} r^h$                                    // restriction
6    $e^{2h} \leftarrow 0$ 
7    $e^{2h} \leftarrow V\text{-cycle}^{2h}(e^{2h}, r^{2h})$ 
8    $e^h \leftarrow I_{2h}^h e^{2h}$                                      // interpolation
9    $v^h \leftarrow v^h + e^h$                                        // correction
10 Relax  $v^h$  for  $n_3$  iterations                                // post smoothing
11 return  $v^h$ 
```

MULTIGRID V-CYCLE

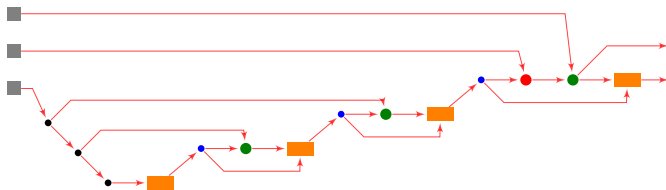


(a) V-cycle



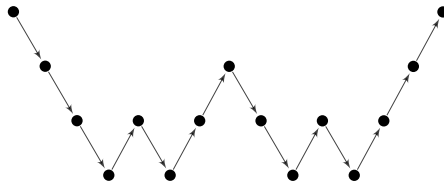
(b) V-cycle: complete DAG

NAS MG V-CYCLE

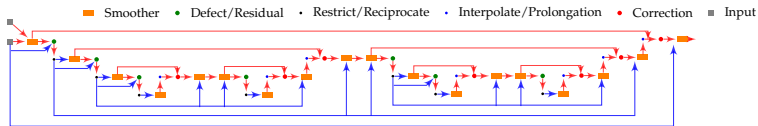


(c) NAS-PB MG V-cycle

MULTIGRID W-CYCLE



(d) W-cycle



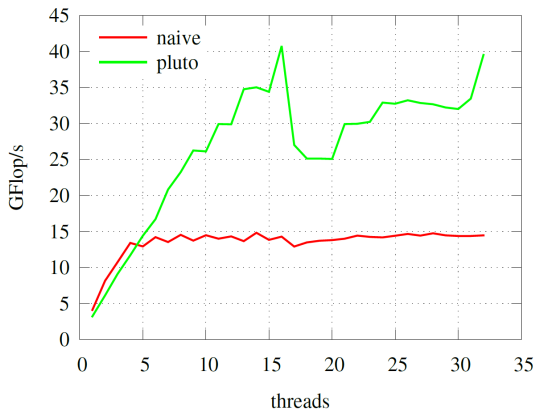
(e) W-cycle: complete DAG

Figure: DAG representation of (a) V-cycle and (b) W-cycle

Strongly recommend reading:

- P. Ghysels and W. Vanroose, Modeling the performance of geometric multigrid on many-core computer architectures, SIAM J. Scientific Computing (2015).
- W. Vanroose, P. Ghysels, D. Roose, and K. Meerbergen, Hiding global communication latency and increasing the arithmetic intensity in extreme-scale Krylov solvers, Position Paper at DOE/ASCR workshop on Applied Mathematics Research for Exascale Computing. Aug 2013.

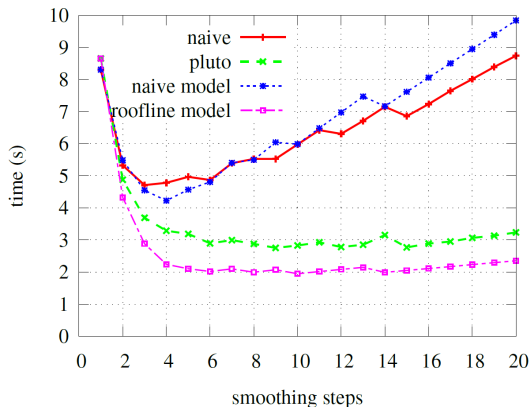
GMG: SMOOTHER SCALING



Scalability of 10 iterations of the Jacobi smoother on an 8000^2 domain on a 16-core Intel Sandy Bridge

Source: Ghysels (LBNL) and Vanroose (University of Antwerp) SIAM J. Scientific Computing 2015

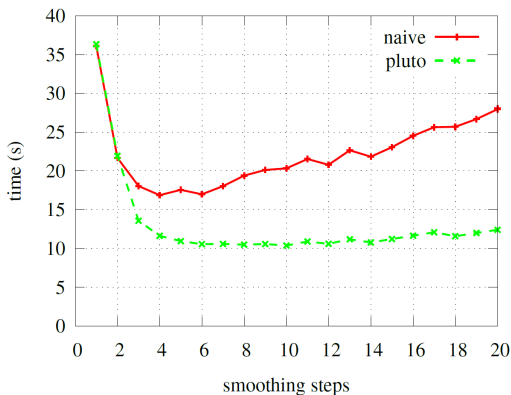
GMG: EXECUTION TIME (2-D)



Timings for a full solve on a 8191^2 domain using V-cycles with a relative stopping tolerance 10^{-12}

Source: Ghysels and Vanroose 2015

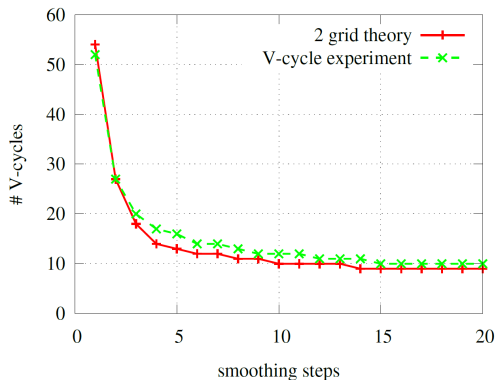
GMG: EXECUTION TIME (3-D)



Timings for a full solve on a 511^3 domain using V-cycles with a relative stopping tolerance 10^{-12} on a dual socket Sandy Bridge machine for a 3D domain

Source: Ghysels and Vanroose 2015

GMG: CONVERGENCE FOR SMOOTHING STEPS

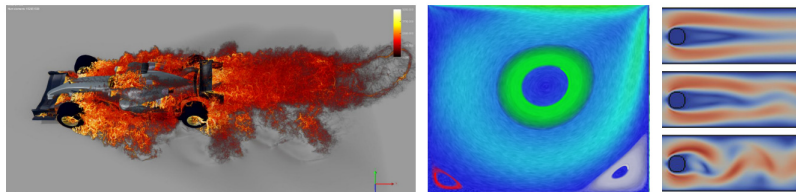


The corresponding number of V-cycles required to reach a 10^{-12} relative stopping criterion for both two-grid and multigrid.

Source: Ghysels and Vanroose 2015

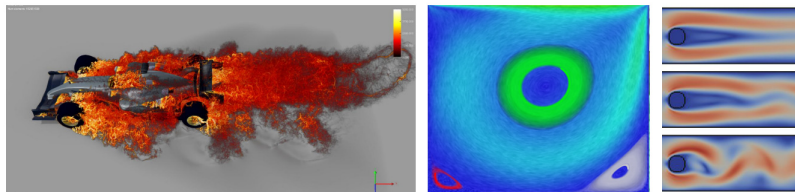
- 1 Pluto/Pluto+
 - Affine Transformations
 - Tiling
- 2 Case Studies
 - Solving Partial Differential Equations
 - Lattice Boltzmann Method
 - Image Processing Pipelines
- 3 Conclusions

LBM: INTRODUCTION



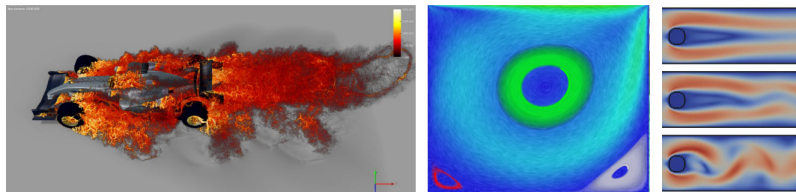
- Lattice-Boltzmann method (LBM) is used for simulation of complex fluid flows in Computational Fluid Dynamics (CFD)

LBM: INTRODUCTION



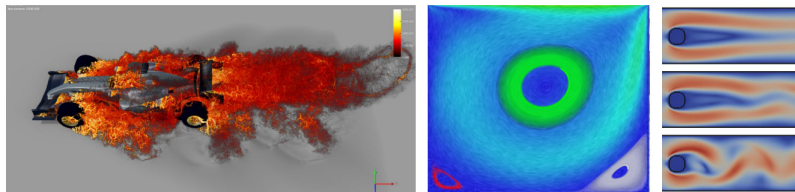
- Lattice-Boltzmann method (LBM) is used for simulation of complex fluid flows in Computational Fluid Dynamics (CFD)
- The simplicity of formulation and its versatility explain the rapid expansion of LBM to applications in complex and multiscale flows

LBM: INTRODUCTION



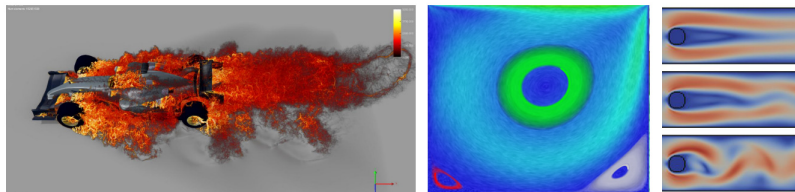
- Lattice–Boltzmann method (LBM) is used for simulation of complex fluid flows in Computational Fluid Dynamics (CFD)
- The simplicity of formulation and its versatility explain the rapid expansion of LBM to applications in complex and multiscale flows
- Particularly suited for parallel and high performance implementations

LBM: INTRODUCTION



- Lattice–Boltzmann method (LBM) is used for simulation of complex fluid flows in Computational Fluid Dynamics (CFD)
- The simplicity of formulation and its versatility explain the rapid expansion of LBM to applications in complex and multiscale flows
- Particularly suited for parallel and high performance implementations
- In spite of tremendous advances in its application, several fundamental opportunities for optimization remain

LBM: INTRODUCTION



- Lattice-Boltzmann method (LBM) is used for simulation of complex fluid flows in Computational Fluid Dynamics (CFD)
- The simplicity of formulation and its versatility explain the rapid expansion of LBM to applications in complex and multiscale flows
- Particularly suited for parallel and high performance implementations
- In spite of tremendous advances in its application, several fundamental opportunities for optimization remain
- We explore one such opportunity through this work

- Fluid flows are modelled as hypothetical particles
 - moving in a lattice domain (discretized space)
 - with different lattice velocities (discretized momentum)
 - over different time steps (discretized time)

- Fluid flows are modelled as hypothetical particles
 - moving in a lattice domain (discretized space)
 - with different lattice velocities (discretized momentum)
 - over different time steps (discretized time)
- Solves the discrete Boltzmann equation for the particle distribution function (a probability density function)

LBM - LATTICE ARRANGEMENTS

- Lattice arrangements are represented as $DmQn$
 - m is the space dimensionality of the lattice
 - n is the number of PDFs (or speeds) involved

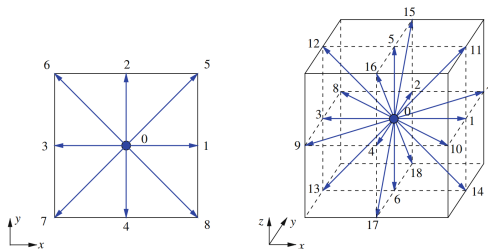


Figure: $D2Q9$ (left) & $D3Q19$ (right) lattice arrangements

LBM - LATTICE ARRANGEMENTS

- Lattice arrangements are represented as $DmQn$
 - m is the space dimensionality of the lattice
 - n is the number of PDFs (or speeds) involved
- Choice of lattice affects precision and duration of simulation

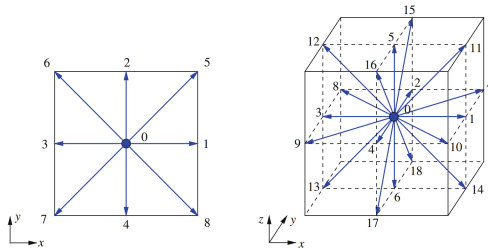


Figure: D2Q9 (left) & D3Q19 (right) lattice arrangements

- The discretized form of Lattice–Boltzmann Equation forms the basis of all LBM models

$$f_i(\mathbf{x} + \mathbf{c}_i \Delta t, t + \Delta t) = f_i(\mathbf{x}, t) + \Omega_i(f_i(\mathbf{x}, t)), \quad i = 1, \dots, n. \quad (1)$$

- The discretized form of Lattice–Boltzmann Equation forms the basis of all LBM models

$$f_i(\mathbf{x} + \mathbf{c}_i \Delta t, t + \Delta t) = f_i(\mathbf{x}, t) + \Omega_i(f_i(\mathbf{x}, t)), \quad i = 1, \dots, n. \quad (1)$$

- Eqn. 1 is solved in two steps, the **collision** step (Eqn. 2) & the **advection** step (Eqn. 3)

$$f_i^*(\mathbf{x}, t + \Delta t) = f_i(\mathbf{x}, t) + \Omega_i(f_i(\mathbf{x}, t)) \quad (2)$$

$$f_i(\mathbf{x} + \mathbf{c}_i \Delta t, t + \Delta t) = f_i^*(\mathbf{x}, t + \Delta t) \quad (3)$$

LBM - IMPLEMENTATION STRATEGIES

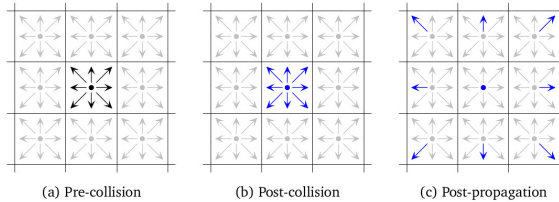


Figure: Push Scheme

LBM - IMPLEMENTATION STRATEGIES

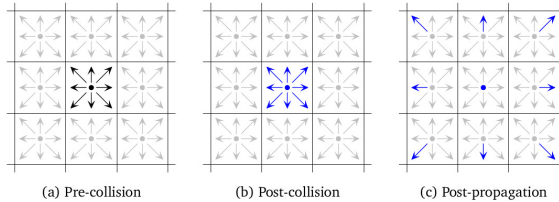


Figure: Push Scheme

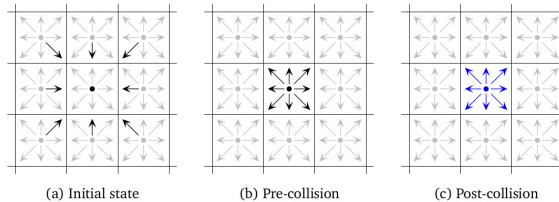


Figure: Pull Scheme

Is it possible to optimize LBM using
time tiling?

TIME TILING LBM COMPUTATIONS

- LBM can be written using storage for either one grid or two grids

TIME TILING LBM COMPUTATIONS

- LBM can be written using storage for either one grid or two grids
- One grid \Rightarrow Separate collision and advection
- Fused Collision + Advection \Rightarrow Two grids

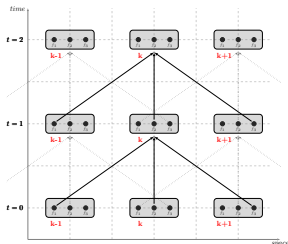


Figure: 1D LBM with single grid

TIME TILING LBM COMPUTATIONS

- LBM can be written using storage for either one grid or two grids
- One grid \Rightarrow Separate collision and advection
- Fused Collision + Advection \Rightarrow Two grids

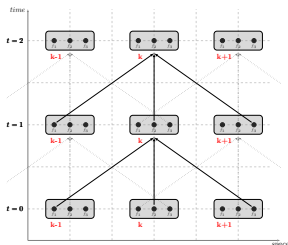


Figure: 1D LBM with single grid

- Not possible to “time tile” LBM with single grid

TIME TILING LBM COMPUTATIONS

- LBM can be written using storage for either one grid or two grids
- One grid \Rightarrow Separate collision and advection
- Fused Collision + Advection \Rightarrow Two grids

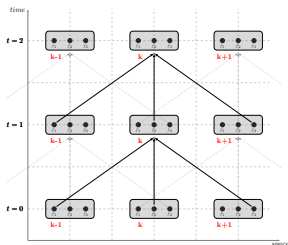


Figure: 1D LBM with single grid

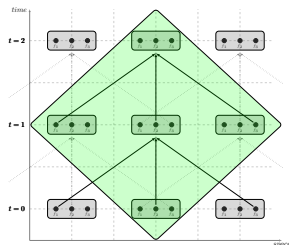


Figure: Pull scheme on a 1D LBM

- Not possible to “time tile” LBM with single grid

TIME TILING LBM COMPUTATIONS

- LBM can be written using storage for either one grid or two grids
- One grid \Rightarrow Separate collision and advection
- Fused Collision + Advection \Rightarrow Two grids

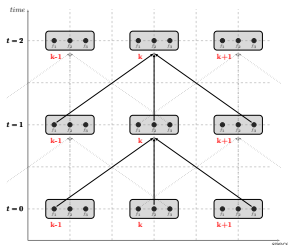


Figure: 1D LBM with single grid

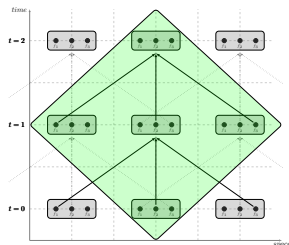


Figure: Pull scheme on a 1D LBM

- Not possible to “time tile” LBM with single grid
- Time tiling is possible with two grids

- We utilize a fused version of the LBM kernel
 - No explicit *advection* phase
 - 2 data grids with a pull scheme to read data
 - Array-of-Structures (AoS) layout for data

INPUT TO POLYHEDRAL TILER

```
#pragma scop
for (t = 0; t < _nTimesteps; t++)
  for (y = 2; y < _nY; y++)
    for (x = 1; x < _nX; x++)
      lbm_kernel(grid[t % 2][y][x][C],
                grid[t % 2][y - 1][x + 0][N],
                grid[t % 2][y + 1][x + 0][S],
                grid[t % 2][y + 0][x - 1][E],
                grid[t % 2][y + 0][x + 1][W],
                grid[t % 2][y - 1][x - 1][NE],
                grid[t % 2][y - 1][x + 1][NW],
                grid[t % 2][y + 1][x - 1][SE],
                grid[t % 2][y + 1][x + 1][SW],

                &grid[(t + 1) % 2][y][x][C],
                &grid[(t + 1) % 2][y][x][N],
                &grid[(t + 1) % 2][y][x][S],
                &grid[(t + 1) % 2][y][x][E],
                &grid[(t + 1) % 2][y][x][W],
                &grid[(t + 1) % 2][y][x][NE],
                &grid[(t + 1) % 2][y][x][NW],
                &grid[(t + 1) % 2][y][x][SE],
                &grid[(t + 1) % 2][y][x][SW], t, y, x);
```

- Use the **PET polyhedral frontend** [Verdoolaeghe and Grosser 2012]: the LBM collision is treated as a blackbox (abstracted as a single function)
- Dependence structure is now similar to “toy time-iterated stencils”
- **All time tiling strategies can now be applied!**

Intel Xeon E5-2680 (SandyBridge)	
Clock	2.7 GHz
Cores / socket	8
Total cores	16
L1 cache / core	32 KB
L2 cache / core	512 KB
L3 cache / socket	20 MB
Peak GFLOPs	172.8
Compiler	Intel C compiler (icc) 14.0.1
Compiler flags	-O3 -xHost -ipo -fno-alias -fno-falias -restrict -fp-model precise -fast-transcendentals
Linux kernel	3.8.0-38

Table: Architecture details

- We compare the performance of our framework on 7 benchmarks against
 - Palabos - an open-source CFD solver based on LBM
 - Compiler auto-parallelization [*icc-auto-par*]
 - Naive manual parallelization using OpenMP [*icc-omp-par*]

- Lid Driven Cavity - d2q9, d3q19 and d3q27

- Lid Driven Cavity - d2q9, d3q19 and d3q27
- SPEC LBM [*470.lbm* from SPEC2006] - d3q19

- Lid Driven Cavity - d2q9, d3q19 and d3q27
- SPEC LBM [*470.lbm* from SPEC2006] - d3q19
- Poiseuille Flow - d2q9

- Lid Driven Cavity - d2q9, d3q19 and d3q27
- SPEC LBM [*470.lbm* from SPEC2006] - d3q19
- Poiseuille Flow - d2q9
- Flow Past Cylinder - d2q9

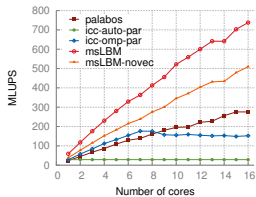
- Lid Driven Cavity - d2q9, d3q19 and d3q27
- SPEC LBM [*470.lbm* from SPEC2006] - d3q19
- Poiseuille Flow - d2q9
- Flow Past Cylinder - d2q9
- MRT - GLBM - d2q9

- Lid Driven Cavity - d2q9, d3q19 and d3q27
- SPEC LBM [*470.lbm* from SPEC2006] - d3q19
- Poiseuille Flow - d2q9
- Flow Past Cylinder - d2q9
- MRT - GLBM - d2q9

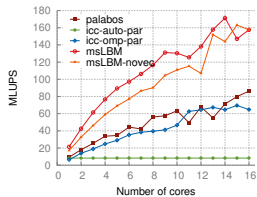
Performance Metrics

- MLUPS - Million Lattice site Updates Per Second
- MEUPS - Million Element Updates Per Second

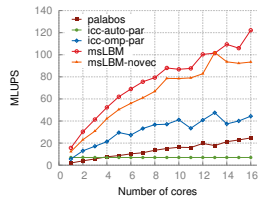
PERFORMANCE - LDC



(a) ldc-d2q9

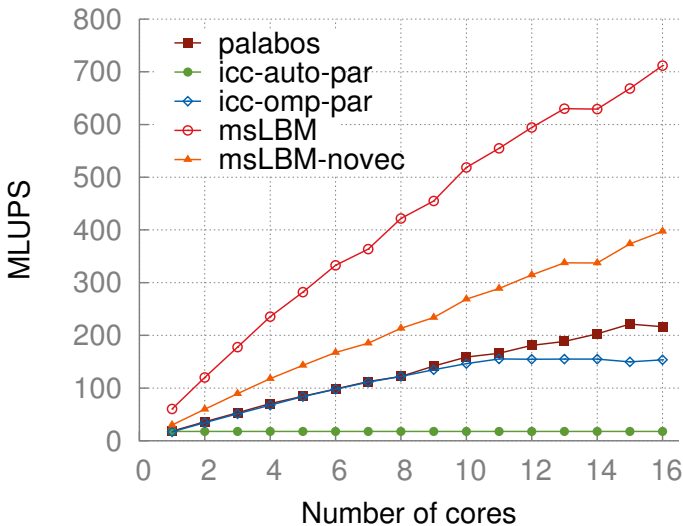


(b) ldc-d3q19



(c) ldc-d3q27

PERFORMANCE - MRT (D2Q9)



ROOFLINE PERFORMANCE MODEL - CONTD.

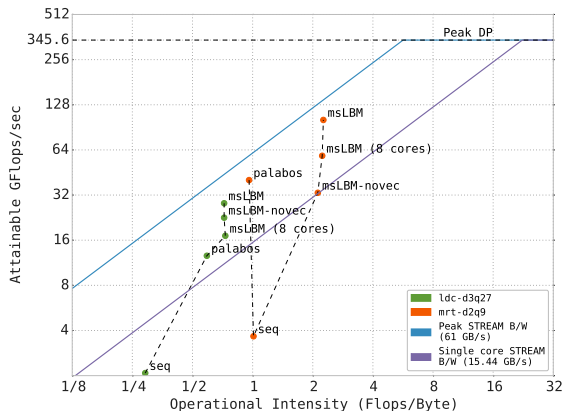


Figure: Roofline model for mrt-d2q9 & ldc-d3q27

ROOFLINE PERFORMANCE MODEL - CONTD.

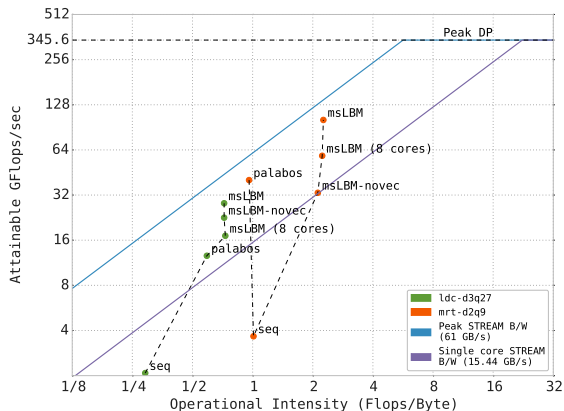


Figure: Roofline model for mrt-d2q9 & ldc-d3q27

- msLBM obtains further improvement over Palabos in both operational intensity and peak achievable performance

- 1 Pluto/Pluto+
 - Affine Transformations
 - Tiling
- 2 Case Studies
 - Solving Partial Differential Equations
 - Lattice Boltzmann Method
 - Image Processing Pipelines
- 3 Conclusions

PolyMage

<http://mcl.csa.iisc.ernet.in/polymage.html>

A DSL and Compiler for Automatic Parallelization and
Optimization of Image Processing Pipelines

Graphs of interconnected processing stages

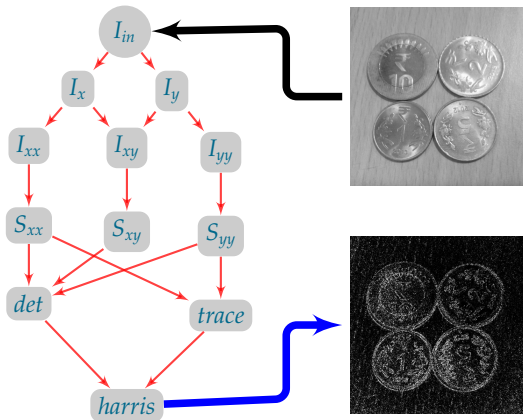


Figure: Harris corner detection

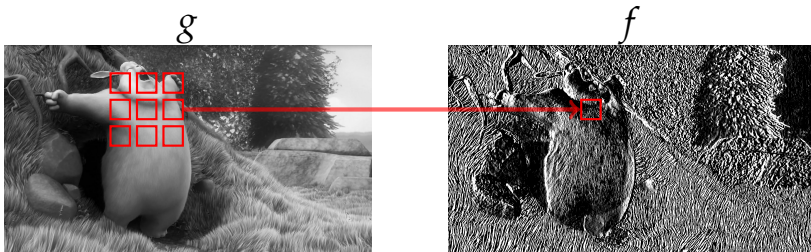
COMPUTATION PATTERNS



Point-wise

$$f(x, y) = w_r \cdot g(x, y, \bullet) + w_g \cdot g(x, y, \bullet) + w_b \cdot g(x, y, \bullet)$$

COMPUTATION PATTERNS



Stencil

$$f(x, y) = \sum_{\sigma_x=-1}^{+1} \sum_{\sigma_y=-1}^{+1} g(x + \sigma_x, y + \sigma_y) \cdot w(\sigma_x, \sigma_y)$$

COMPUTATION PATTERNS



Downsample

$$f(x, y) = \sum_{\sigma_x=-1}^{+1} \sum_{\sigma_y=-1}^{+1} g(2x + \sigma_x, 2y + \sigma_y) \cdot w(\sigma_x, \sigma_y)$$

COMPUTATION PATTERNS



$$f(x, y) = \sum_{\sigma_x=-1}^{+1} \sum_{\sigma_y=-1}^{+1} g((x + \sigma_x)/2, (y + \sigma_y)/2) \cdot w(\sigma_x, \sigma_y, x, y)$$

EXAMPLE: PYRAMID BLENDING PIPELINE

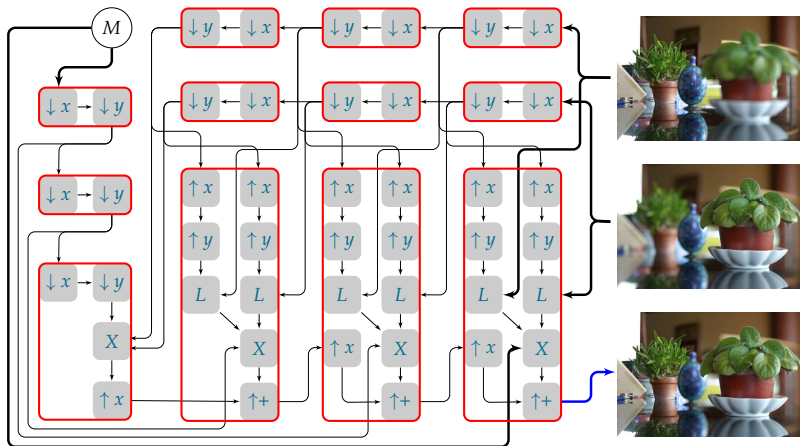


Image courtesy: Kyros Kutulakos

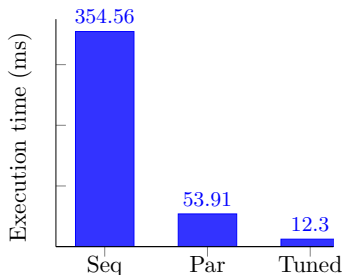
WHERE ARE IMAGE PROCESSING PIPELINES USED?

- On images uploaded to social networks like Facebook, Google+
- On all camera-enabled devices
- Everyday workloads from data center to mobile device scales
- **Computational photography, computer vision, medical imaging, ...**

Google+ Auto Enhance



NAIVE VS OPTIMIZED IMPLEMENTATION

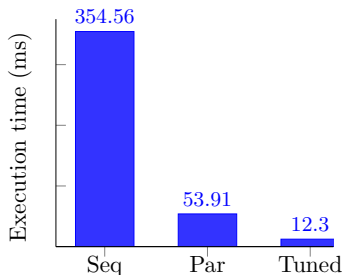


- Naive implementation in C
- Naive parallelization – $7\times$ OpenMP, Vector pragmas (icc)
- Manual optimization – $29\times$ Locality, Parallelism, Vector intrinsics

**Harris corner detection
(16 cores)**

Manually optimizing pipelines is hard

NAIVE VS OPTIMIZED IMPLEMENTATION



- Naive implementation in C
- Naive parallelization – $7\times$ OpenMP, Vector pragmas (icc)
- Manual optimization – $29\times$ Locality, Parallelism, Vector intrinsics

**Harris corner detection
(16 cores)**

**Goal: Performance levels of manual tuning
Without the pain**

OUR APPROACH: POLYMAGE

- **High-level language (DSL embedded in Python)**
 - Allow expressing common patterns intuitively
 - Enables compiler analysis and optimization
- **Automatic Optimizing Code Generator**
 - Uses domain-specific cost models to apply complex combinations of scaling, alignment, **tiling** and **fusion** to optimize for **parallelism** and **locality**

HARRIS CORNER DETECTION

```

R, C = Parameter(Int), Parameter(Int) (*\label{param}*)
I = Image(Float, [R+2, C+2]) (*\label{image}*)

x, y = Variable(), Variable() (*\label{vars}*)
row, col = Interval(0,R+1,1), Interval(0,C+1,1) (*\label{intervals}*)

c = Condition(x,'>=',1) & Condition(x,'<=',R) & (*\label{cond1}*)
    Condition(y,'>=',1) & Condition(y,'<=',C)

cb = Condition(x,'>=',2) & Condition(x,'<=',R-1) & (*\label{cond2}*)
    Condition(y,'>=',2) & Condition(y,'<=',C-1)

Iy = Function(varDom = ([x,y],[row,col]),Float) (*\label{f1}*)
Iy.defn = [ Case(c, Stencil(I(x,y), 1.0/12, (*\label{d1}*)
    [[-1, -2, -1],
     [ 0, 0, 0],
     [ 1, 2, 1]]) )

Ix = Function(varDom = ([x,y],[row,col]),Float) (*\label{f2}*)
Ix.defn = [ Case(c, Stencil(I(x,y), 1.0/12, (*\label{d2}*)
    [[-1, 0, 1],
     [-2, 0, 2],
     [-1, 0, 1]]) )

Ixx = Function(varDom = ([x,y],[row,col]),Float) (*\label{f3}*)
Ixx.defn = [ Case(c, Ix(x,y) * Ix(x,y)) (*\label{d3}*)

Iyy = Function(varDom = ([x,y],[row,col]),Float) (*\label{f4}*)
Iyy.defn = [ Case(c, Iy(x,y) * Iy(x,y)) (*\label{d4}*)

Ixy = Function(varDom = ([x,y],[row,col]),Float) (*\label{f5}*)
Ixy.defn = [ Case(c, Ix(x,y) * Iy(x,y)) (*\label{d5}*)

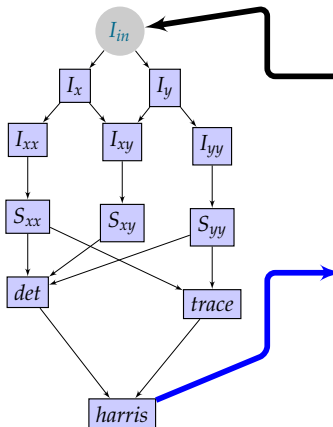
Sxx = Function(varDom = ([x,y],[row,col]),Float) (*\label{f6}*)
Syy = Function(varDom = ([x,y],[row,col]),Float) (*\label{f7}*)
Sxy = Function(varDom = ([x,y],[row,col]),Float) (*\label{f8}*)
for pair in [(Sxx, Ixx), (Syy, Iyy), (Sxy, Ixy)]: (*\label{meta}*)
    pair[0].defn = [ Case(cb, Stencil(pair[1], 1, (*\label{d6}*)
        [[1, 1, 1],
         [1, 1, 1],
         [1, 1, 1]]) )

det = Function(varDom = ([x,y],[row,col]),Float) (*\label{f9}*)
d = Sxx(x,y) * Syy(x,y) - Sxy(x,y) * Sxy(x,y)
det.defn = [ Case(cb, d) ] (*\label{d7}*)

trace = Function(varDom = ([x,y],[row,col]),Float) (*\label{f10}*)
trace.defn = [ Case(cb, Sxx(x,y) + Syy(x,y)) ] (*\label{d8}*)

harris = Function(varDom = ([x,y],[row,col]),Float) (*\label{f11}*)
coarsity = det(x,y) - .04 * trace(x,y) * trace(x,y) (*\label{d9}*)
harris.defn = [ Case(cb, coarsity) ]

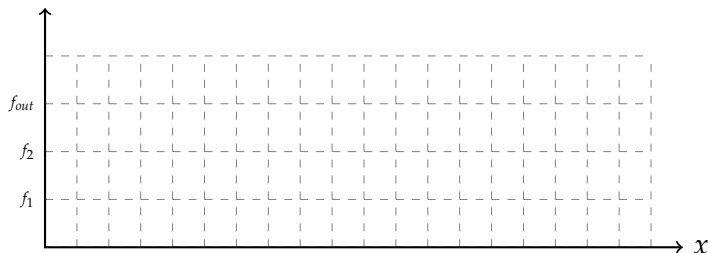
```



OUR APPROACH: POLYMAGE

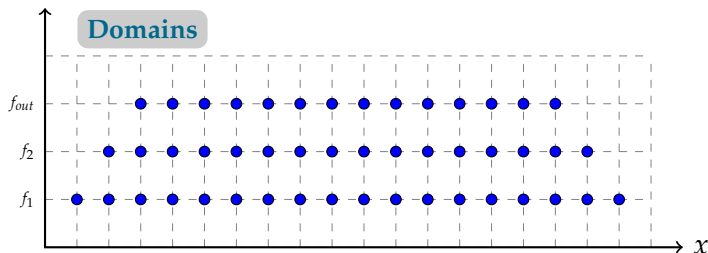
- High-level language (DSL embedded in Python)
 - Allow expressing common patterns intuitively
 - Enables compiler analysis and optimization
- **Automatic Optimizing Code Generator**
 - Uses domain-specific cost models to apply complex combinations of scaling, alignment, **tiling** and **fusion** to optimize for **parallelism** and **locality**

POLYHEDRAL REPRESENTATION



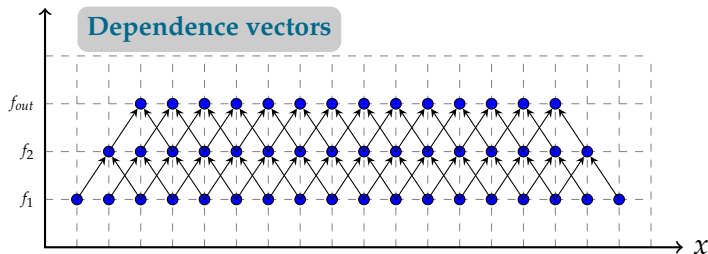
```
x = Variable()
f_in = Image(Float, [18])
f1 = Function(varDom = ([x], [Interval(0, 17, 1)]), Float)
f1.defn = [ f_in(x) + 1 ]
f2 = Function(varDom = ([x], [Interval(1, 16, 1)]), Float)
f2.defn = [ f1(x-1) + f1(x+1) ]
f_out = Function(varDom = ([x], [Interval(2, 15, 1)]), Float)
f_out.defn = [ f2(x-1) + f2(x+1) ]
```


POLYHEDRAL REPRESENTATION



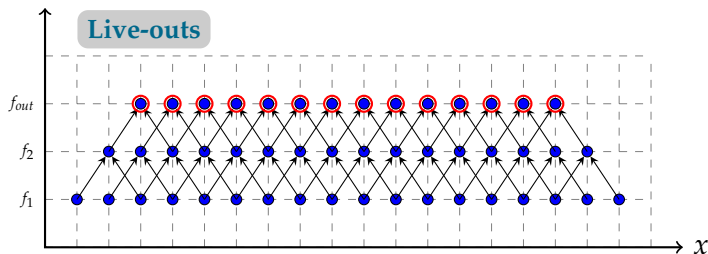
```
x = Variable()
fin = Image(Float, [18])
f1 = Function(varDom = ([x], [Interval(0, 17, 1)]), Float)
f1.defn = [ fin(x) + 1 ]
f2 = Function(varDom = ([x], [Interval(1, 16, 1)]), Float)
f2.defn = [ f1(x-1) + f1(x+1) ]
fout = Function(varDom = ([x], [Interval(2, 15, 1)]), Float)
fout.defn = [ f2(x-1) + f2(x+1) ]
```

POLYHEDRAL REPRESENTATION



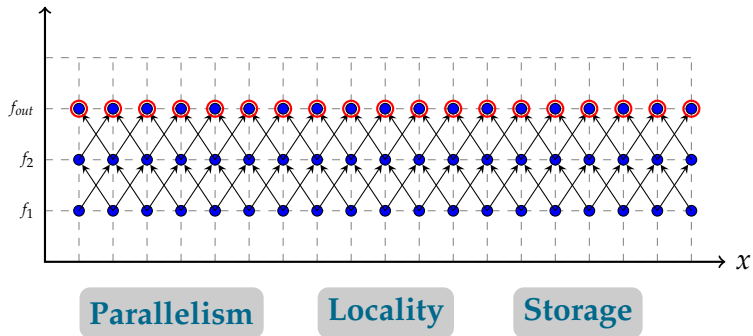
Function	Dependence Vectors
$f_{out}(x) = f_2(x - 1) \cdot f_2(x + 1)$	$(1, 1), (1, -1)$
$f_2(x) = f_1(x - 1) + f_1(x + 1)$	$(1, 1), (1, -1)$
$f_1(x) = f_{in}(x)$	

POLYHEDRAL REPRESENTATION

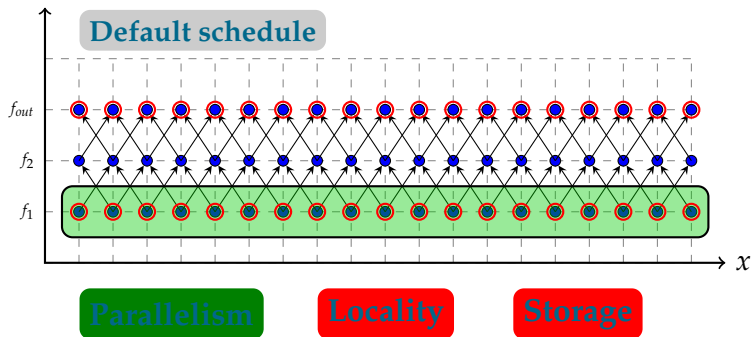


Function	Dependence Vectors
$f_{out}(x) = f_2(x-1) \cdot f_2(x+1)$	$(1, 1), (1, -1)$
$f_2(x) = f_1(x-1) + f_1(x+1)$	$(1, 1), (1, -1)$
$f_1(x) = f_{in}(x)$	

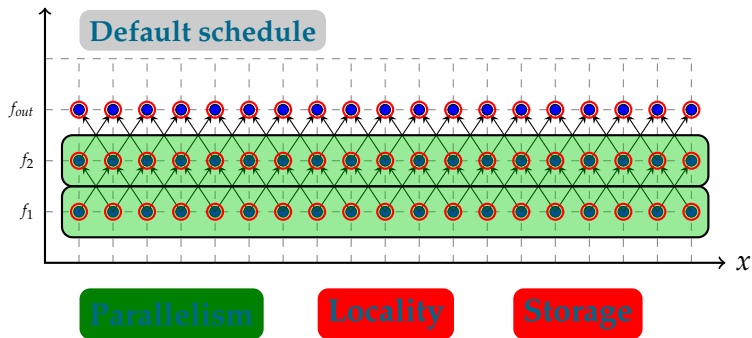
SCHEDULING CRITERIA



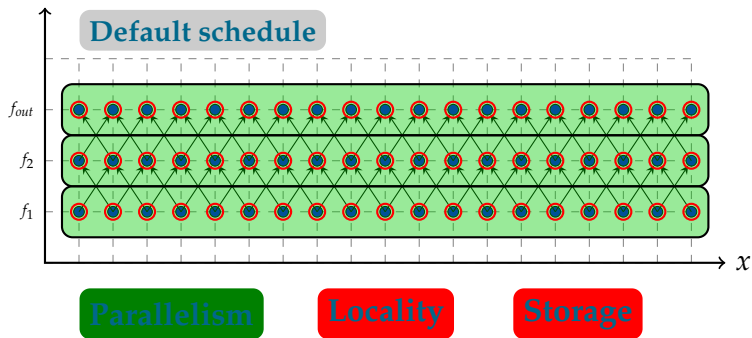
SCHEDULING CRITERIA



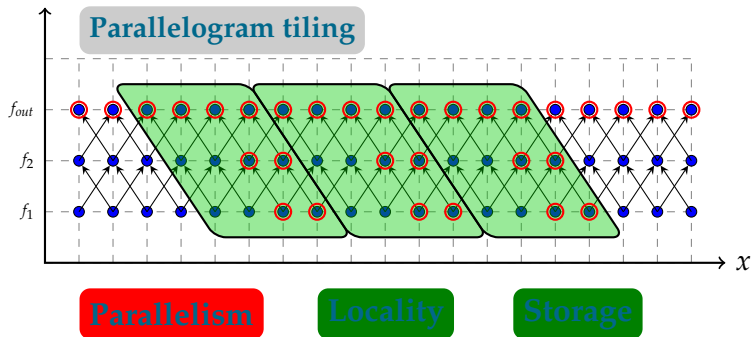
SCHEDULING CRITERIA



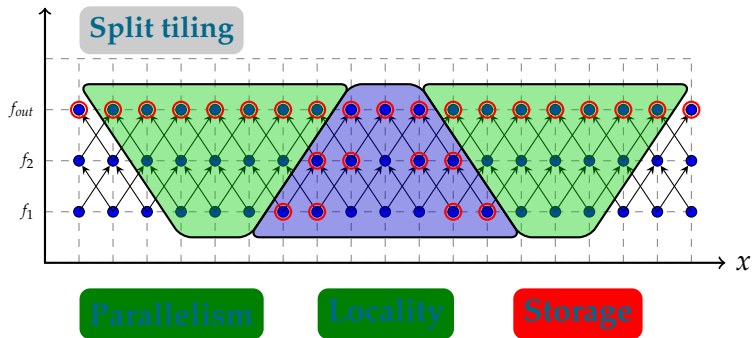
SCHEDULING CRITERIA



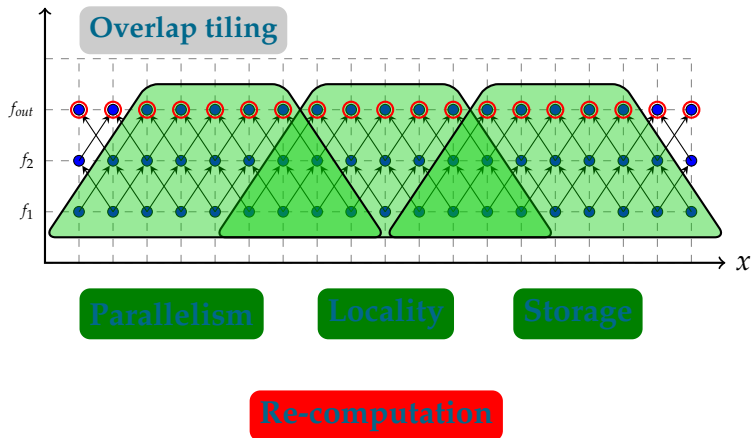
SCHEDULING CRITERIA



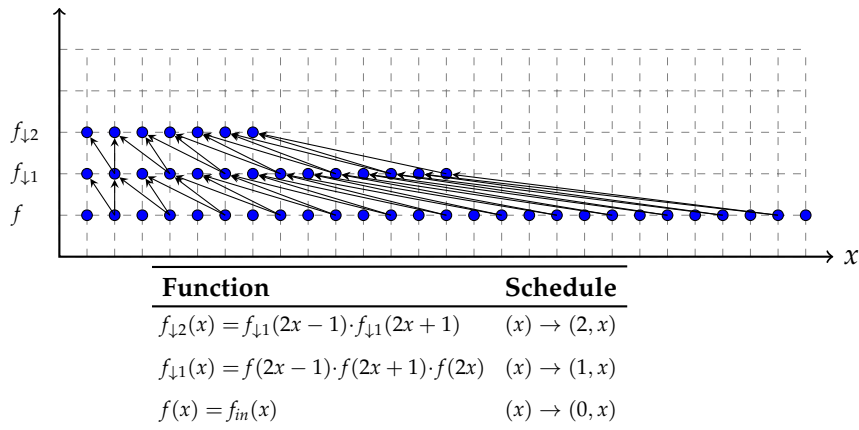
SCHEDULING CRITERIA



SCHEDULING CRITERIA

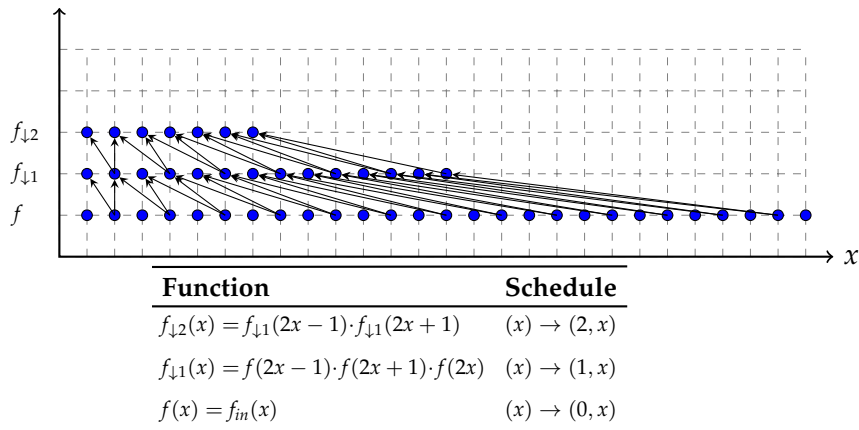


OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



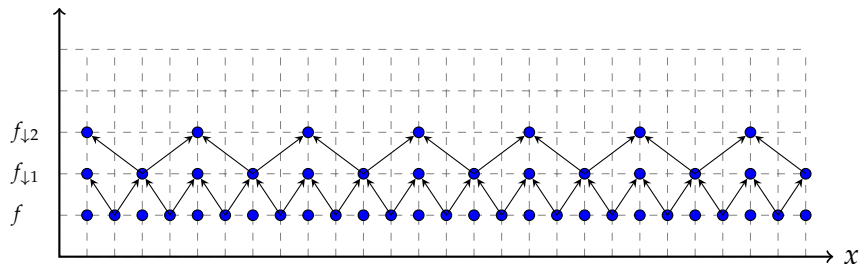
- Prior approaches for overlapped tiling only consider homogeneous time-iterated stencils

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



- Cannot have a fixed tile shape when dependence vectors are non-constant

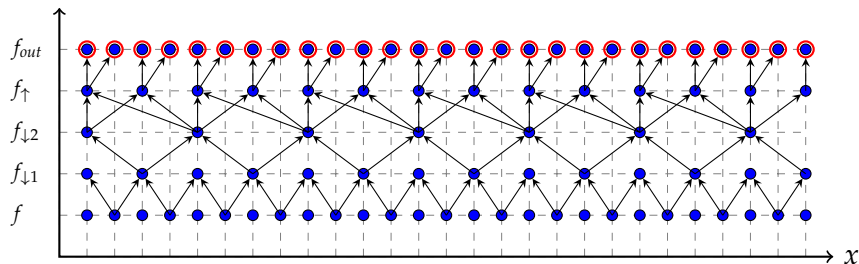
OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



Function	Schedule
$f_{\downarrow 2}(x) = f_{\downarrow 1}(2x - 1) \cdot f_{\downarrow 1}(2x + 1)$	$(x) \rightarrow (2, 4x)$
$f_{\downarrow 1}(x) = f(2x - 1) \cdot f(2x + 1) \cdot f(2x)$	$(x) \rightarrow (1, 2x)$
$f(x) = f_{in}(x)$	$(x) \rightarrow (0, x)$

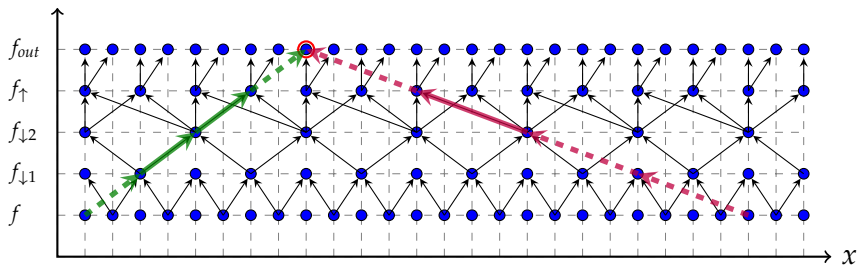
- Scaling and aligning the schedules

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



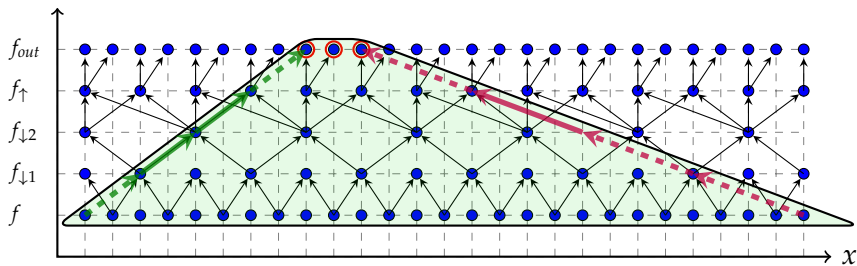
Function	Schedule
$f_{out}(x) = f_{\uparrow}(x/2)$	$(x) \rightarrow (4, x)$
$f_{\uparrow}(x) = f_{\downarrow 2}(x/2) \cdot f_{\downarrow 2}(x/2 + 1)$	$(x) \rightarrow (3, 2x)$
$f_{\downarrow 2}(x) = f_{\downarrow 1}(2x - 1) \cdot f_{\downarrow 1}(2x + 1)$	$(x) \rightarrow (2, 4x)$
$f_{\downarrow 1}(x) = f(2x - 1) \cdot f(2x + 1) \cdot f(2x)$	$(x) \rightarrow (1, 2x)$
$f(x) = f_{in}(x)$	$(x) \rightarrow (0, x)$

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



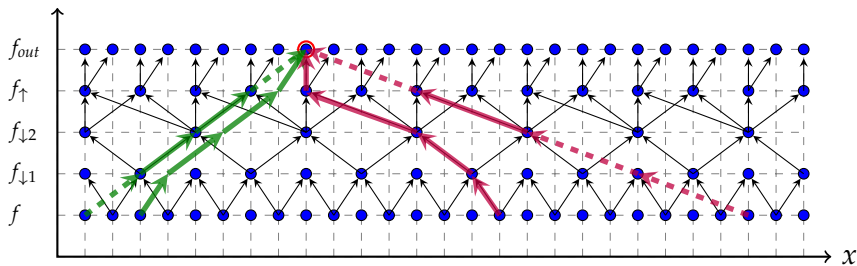
- **Determining tile shape**
Conservative vs precise bounding faces

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



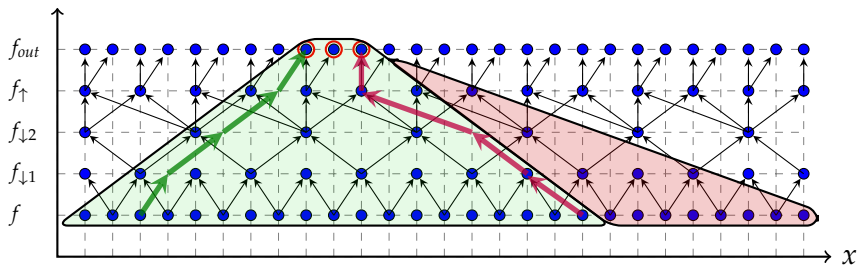
- **Determining tile shape**
Conservative vs precise bounding faces

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



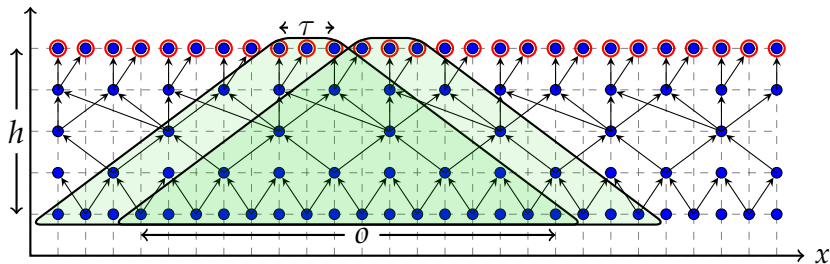
- **Determining tile shape**
Conservative vs precise bounding faces

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



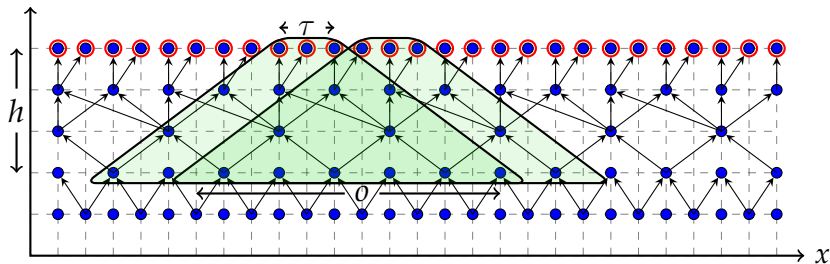
- Significant reduction in redundant computation

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



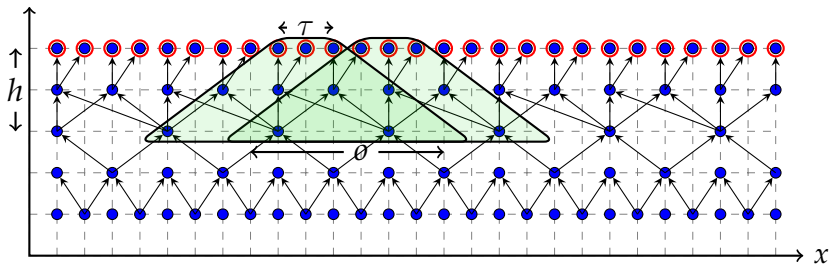
- Tile size τ , Overlap O , Height h
Trade-off between fusion height and overlap

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



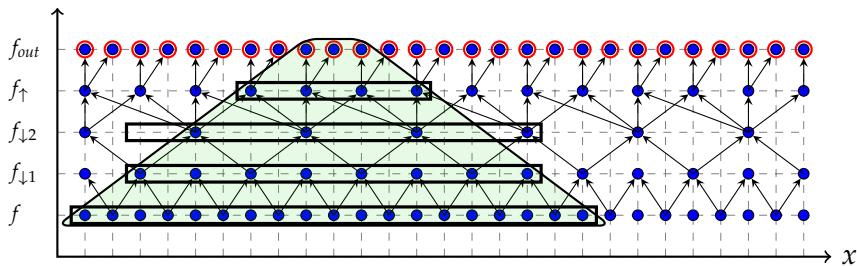
- Tile size τ , Overlap O , Height h
Trade-off between fusion height and overlap

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



- Tile size τ , Overlap O , Height h
Trade-off between fusion height and overlap

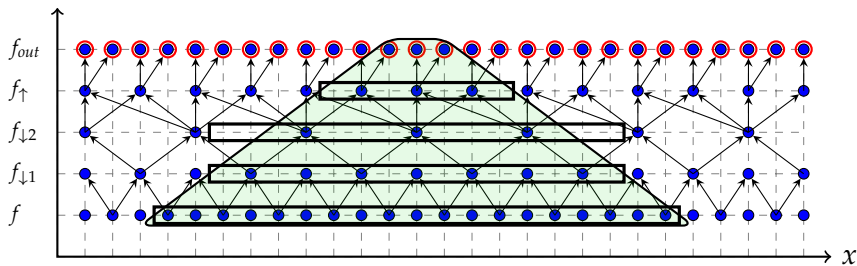
OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



Scratch pads

- Reduction in intermediate storage
- Better locality and reuse
- Privatized for each thread

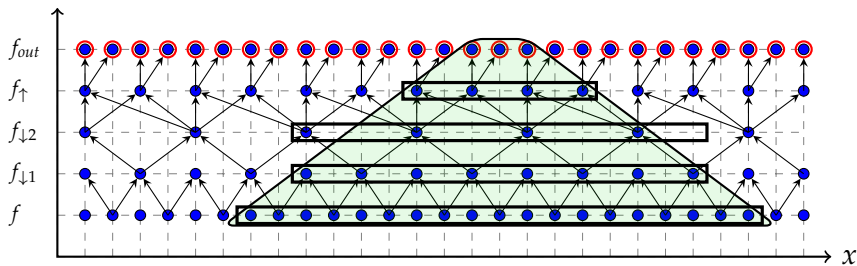
OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



Scratch pads

- Reduction in intermediate storage
- Better locality and reuse
- Privatized for each thread

OVERLAPPED TILING FOR HETEROGENEOUS FUNCTIONS



Scratch pads

- Reduction in intermediate storage
- Better locality and reuse
- Privatized for each thread

Seven representative benchmarks of varying structure and complexity

Benchmark	Stages	Lines	Image size
Unsharp Mask	4	16	$2048 \times 2048 \times 3$
Bilateral Grid	7	43	2560×1536
Harris Corner	11	43	6400×6400
Camera Pipeline	32	86	2528×1920
Pyramid Blending	44	71	$2048 \times 2048 \times 3$
Multiscale Interpolate	49	41	$2560 \times 1536 \times 3$
Local Laplacian	99	107	$2560 \times 1536 \times 3$

OUTLINE

- 1 Pluto/Pluto+
 - Affine Transformations
 - Tiling
- 2 Case Studies
 - Solving Partial Differential Equations
 - Lattice Boltzmann Method
 - Image Processing Pipelines
- 3 Conclusions

- Interesting to see how numerical techniques can be chosen and designed around modern parallel architectures and optimization infrastructure

ACKNOWLEDGEMENTS

- *Ravi Teja Mullapudi, Vinay Vasista, Irshad Pananilath, Aravind Acharya, Vinayaka Bandishti* (students at IISc)
- **Google Research**
- *Anand Deshpande, Aniruddha Shet, Bharat Kaul, and Sunil Sherlekar* from **Intel Labs, Bangalore** for discussions

Thank You!