

V. Elango<sup>†</sup>, F. Rastello<sup>§</sup>, L-N Pouchet<sup>†</sup>, J. Ramanujam<sup>‡</sup>, P. Sadayappan<sup>†</sup>  
<sup>†</sup>The Ohio State University      <sup>§</sup>INRIA      <sup>‡</sup>Louisiana State University

### Modeling Data Movement Complexity

```

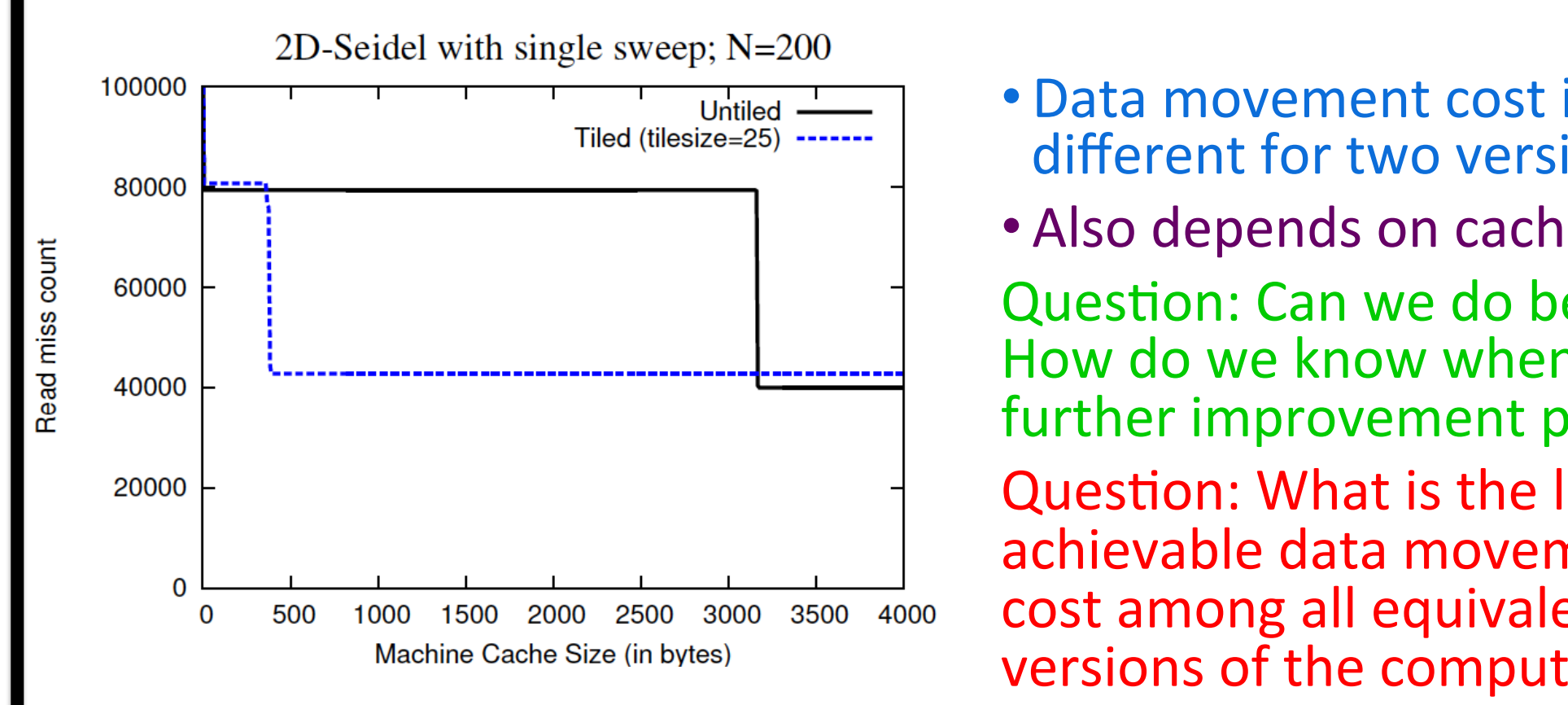
for (i=1; i<N-1; i++)
  for (j=1; j<N-1; j++)
    A[i][j] = A[i][j-1] + A[i-1][j];
  
```

Untiled version  
Comp. complexity:  $(N-1)^2$  Ops

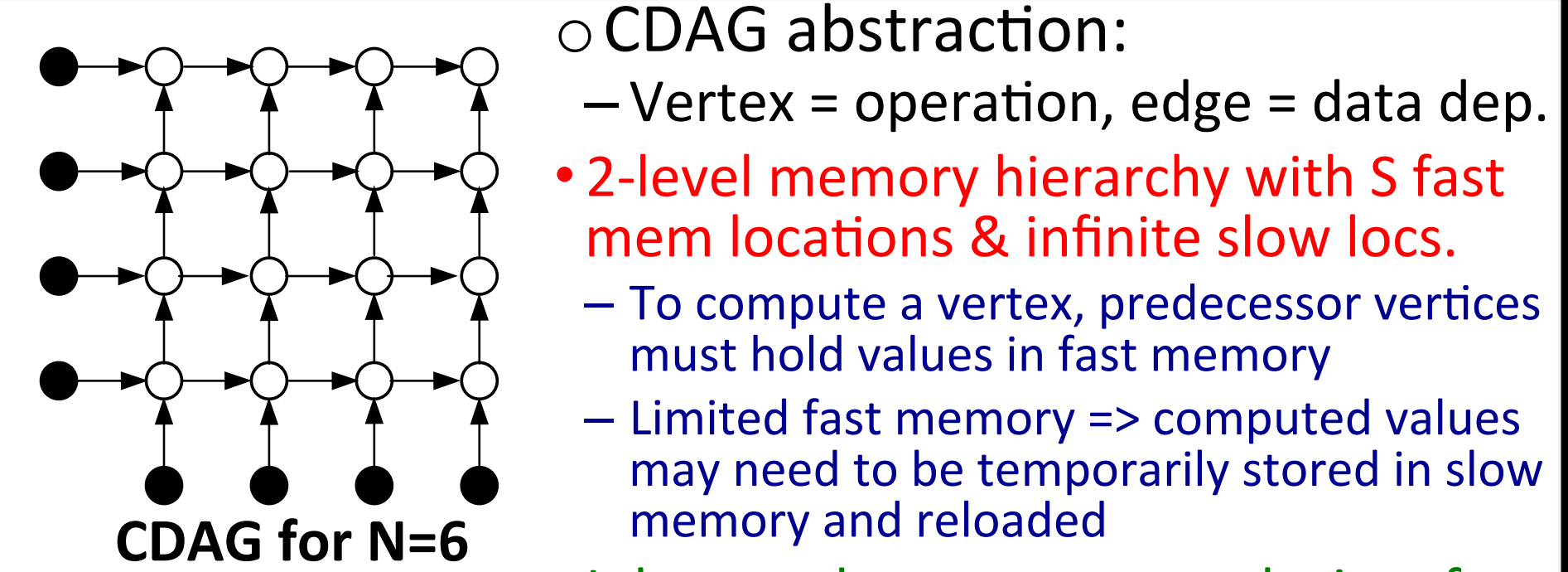
```

for (it = 1; it < N-1; it += B)
  for (j = 1; j < N-1; j += B)
    for (i = it; i < min(it+B, N-1); i++)
      for (k = j; k < min(j+B, N-1); k++)
        A[i][k] = A[i-1][k] + A[i][k-1];
  
```

Tiled Version  
Comp. complexity:  $(N-1)^2$  Ops



• Data movement cost is different for two versions  
 • Also depends on cache size  
 Question: Can we do better? How do we know when no further improvement possible?  
 Question: What is the lowest achievable data movement cost among all equivalent versions of the computation?



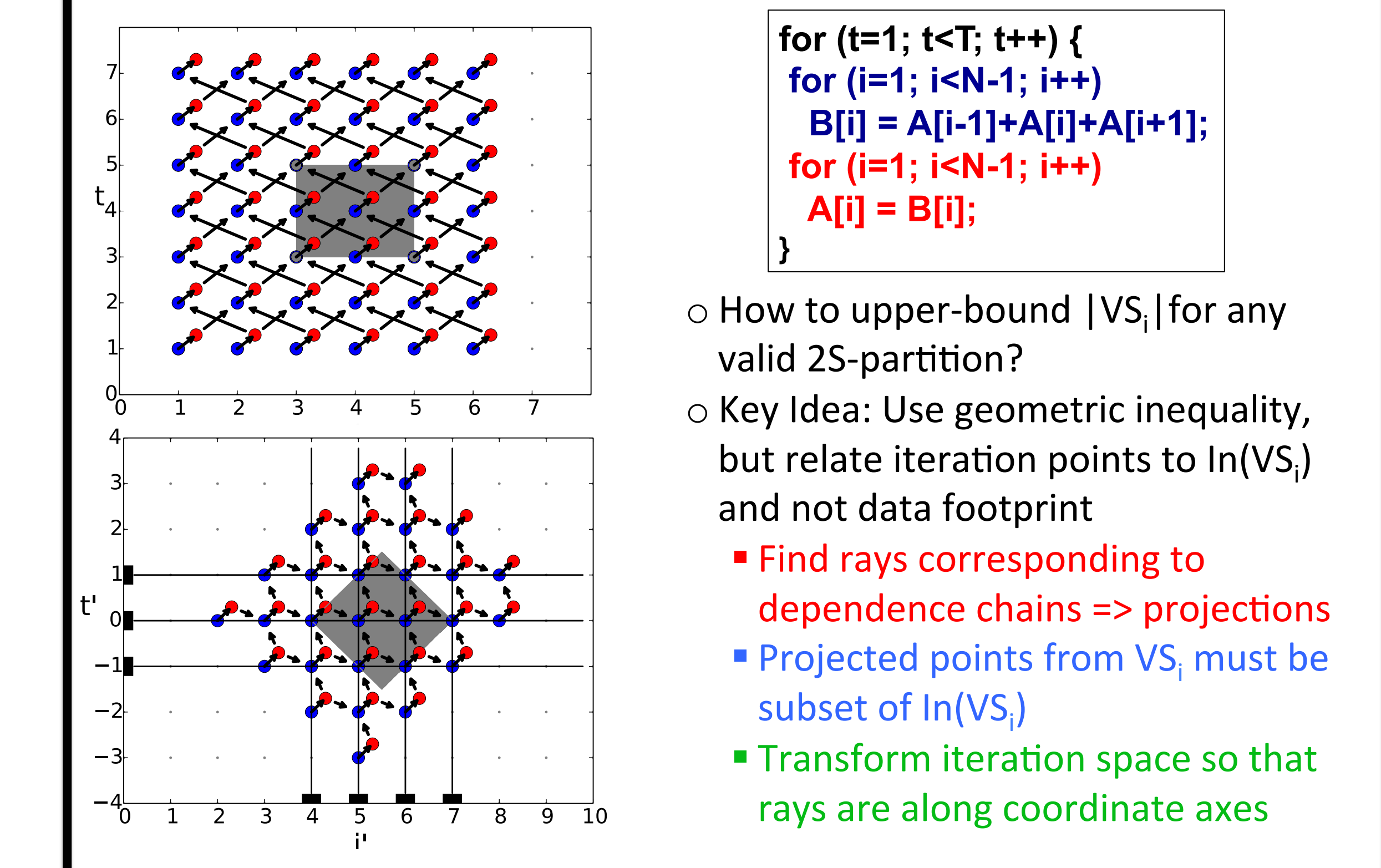
CDAG abstraction:  
 - Vertex = operation, edge = data dep.  
 • 2-level memory hierarchy with  $S$  fast mem locations & infinite slow locs.  
 - To compute a vertex, predecessor vertices must hold values in fast memory  
 - Limited fast memory => computed values may need to be temporarily stored in slow memory and reloaded  
 • Inherent data mvmt. complexity of CDAG: Minimal #loads+#stores among all possible valid schedules

Develop upper bounds on min-cost  
 Minimum possible data movement cost? No known effective solution to problem  
 Develop lower bounds on min-cost

### Prior Work: Data Movement Lower Bounds

- Arbitrary CDAGs:
    - Hong & Kung (1981): strong relation between: 1) Data movement cost for a CDAG schedule, and 2) Number of vertex-sets in "2S-partition" of CDAG
    - Change from reasoning about all valid schedules to all valid 2S-partitions of graph
    - Generalization
    - Manual CDAG-specific reasoning => challenge to automate
  - Linear-Algebra-like algorithms:
    - Irony et al. (2004) and Ballard et al. (2011): Geometric approach based on geometric inequality
    - Christ et al. (2013): Automation, based on generalized geometric HBL inequality (Holder-Brascamp-Lieb)
    - Automated asymptotic parametric lower bound expressions, e.g.,  $O(N^3/\sqrt{S})$  for  $N \times N$  mat-mult
    - Restricted computational model: weakness of bounds or inapplicability
- Our work: Static analysis to automate asymptotic parametric lower bounds analysis of affine codes for CDAG model

### Lower Bounds for CDAGs: Geometric Reasoning



How to upper-bound  $|VS_i|$  for any valid 2S-partition?  
 Key Idea: Use geometric inequality, but relate iteration points to  $\ln(VS_i)$  and not data footprint  
 Find rays corresponding to dependence chains => projections  
 Projected points from  $VS_i$  must be subset of  $\ln(VS_i)$   
 Transform iteration space so that rays are along coordinate axes

### Lower Bounds: Geometric Reasoning with Data Footprints

Loomis-Whitney inequality (2D): bounds #points in a set by product of # projected points on coordinate axes  
 Prior work: Uses Loomis-Whitney inequality & generalization (Holder-Brascamp-Lieb) for lower bounds for linear-algebra-like computations  
 Projections of iteration-space points  $\Leftrightarrow$  Data footprint  
 Geometric inequality: Bound max. #of ops for a given # of data moves

for (i=0; i<N; i++)  
 for (j=0; j<N; j++)  
 if (i < j) force[i] += func(pos[i], pos[j])

2D Loomis-Whitney Inequality  
 $|E| \leq \prod_{j=1}^d |\phi_j(E)|^{1/(d-1)}$   
 $|E| \leq |E_i| * |E_j|$

Divide execution trace into segments with  $S$  load/stores (3 in ex.)  
 Within each segment, #distinct elements of  $pos[] \leq 2S$  (up to  $S$  coming into segment in scratchpad and another  $S$  explicitly loaded)  
 For code example, projection of Stmt(i,j) onto  $i$ -axis maps to data element  $pos[i]$ ; similarly for  $j$ -axis Max. # distinct elts of  $pos[i]$  or  $pos[j]$  read in any segment  $\leq 2S$   
 By Loomis-Whitney, max. # iteration points in any segment,  $|P| \leq 2S * 2S$   
 Min. #segments  $\geq N^2/4S^2$ ; each seg. (but last) has  $S$  load/stores  
 #load/stores  $\geq (N^2/4S^2 - 1) * S = \Omega(N^2/S)$

### CDAG Lower Bounds: Hong/Kung S-Partitioning

$P1 \forall i \neq j, V_i \cap V_j = \emptyset$ , and  $\bigcup_{i=1}^n V_i = V \setminus I$   
 $P2$  there is no cyclic dependence between subsets  
 $P3 \forall i, |\ln(V_i)| \leq S$  S-Partition of CDAG satisfies 4 properties  
 $P4 \forall i, |\text{Out}(V_i)| \leq S$

Any valid schedule using  $S$  registers is associated with a 2S-partition of CDAG  
 Divide trace into segments incurring exactly  $S$  load/stores  
 Ops executed in segment- $i$  form a convex vertex set  $VS_i$   
 $|\ln(VS_i)| \leq 2S$  (up to  $S$  from prev. segment and up to  $S$  new loads)  
 Each segment (except last) has  $S$  loads/stores =>  $S * NS \geq \text{Total I/O} \geq S * (NS - 1)$   
 Reasoning about minimum #vertex sets for any valid 2S-partition => Lower bound on # loads/stores

Parameters:  $N, T$   
 Inputs:  $\ln[N]$ ; Outputs:  $A[N]$   
 for (i=0; i<N; i++)  
 A[i] = ln[i]; S1  
 for (t=0; t<T; t++) {  
 for (i=1; i<N-1; i++)  
 B[i] = A[i-1]+A[i]+A[i+1]; S2  
 for (i=1; i<N-1; i++)  
 A[i] = B[i]; S3

Use ISL to find all "must" data flow dependences  
 Cycles data dep. graph = "rays" in the CDAG  
 Generalized geom. inequality allows different dimensional orthogonal projections  
 Parametric exponents in inequality: sum of weighted ranks of projected subspaces must exceed rank of full iteration space  
 solve a linear program to find optimal weights  
 => asymptotic parametric I/O lower bound for affine program

$|E| \leq \prod_{j=1}^d |\phi_j(E)|^{1/(d-1)} \rightsquigarrow |E| \leq \prod_{j=1}^m |\phi_j(E)|^{s_j}$  s.t.,  $\forall i, 1 \leq \sum_{j=1}^m s_j \delta_{i,j}$   
 where,  $(s_1, \dots, s_m) \in [0, 1]^m$   
 $\phi_j: \mathbb{R}^d \rightarrow \mathbb{R}^d$  are orthogonal projections  
 $\delta_{i,j}: \dim(\phi_j(\text{span}(\vec{e}_i)))$  - where  $\vec{e}_i$ ,  $i$ -th canonical vector.

### Geometric Reasoning with Data Footprints: Limitations

Cannot handle multi-statement programs  
 Computations with very different data mvmt. Rqmts. but same array access footprint => same LB  
 Semantics preserving loop transformations can result in change to lower bound  
 Same access functions => same analysis result LB =  $\Omega(N^2/S)$   
 Semantically equivalent code after loop distribution: but different IO lower bound LB =  $\Omega(N^2)$

for (i=0; i<N; i++)  
 for (j=0; j<N; j++)  
 for (k=0; k<N; k++)  
 C[i][j] += A[i][k]\*B[k][j];

for (i=0; i<N; i++)  
 for (j=0; j<N; j++) {  
 C[i][j] += 1;  
 A[i][j] += 1;  
 B[j][i] += 1;  
 }

for (i,j,k) C[i][j] += 1;  
 for (i,j,k) A[i][j] += 1;  
 for (i,j,k) B[k][j] += 1;

Cannot model effect of data dependences  
 Dependences may impose constraints => higher data movement cost than footprint analysis reveals  
 Example: 1D Jacobi - footprint based geometric analysis cannot derive known LB of  $\Omega(NT/S)$

### Lower Bounds: Research Directions

1) Alternate lower bounds approach (graph min-cut based)  
 2) Composition of lower bounds  
 3) Modeling vertical + horizontal data movement bounds for scalable parallel systems [SPAA '14]

Theory & Models

Tools

1) Automated lower bounds for arbitrary explicit CDAGs  
 2) Automated parametric lower bounds for affine programs [this poster; POPL '15]

Applications

1) Comparative analysis of algorithms via lower bounds  
 2) Assessment of compiler effectiveness  
 3) Algorithm/architecture co-design space exploration [ACM TACO '14, Hipec '15]