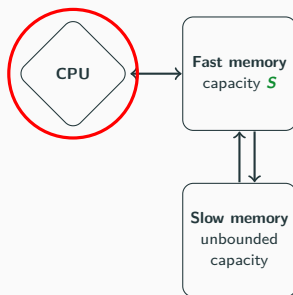


IOOpt: Automatic Derivation of I/O Complexity Bounds for Affine Programs

Auguste Olivry Guillaume Iooss Nicolas Tollenaere
Atanas Rountev P. Sadayappan Fabrice Rastello
June 2021

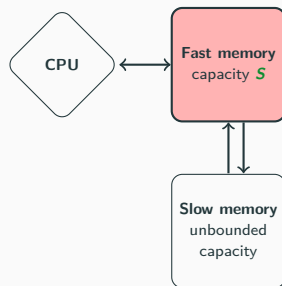
What is I/O complexity?

- Arithmetic complexity = # of operations



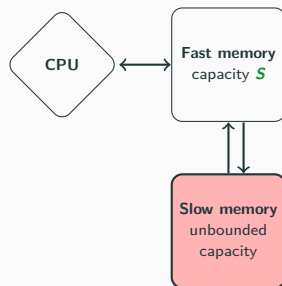
What is I/O complexity?

- Arithmetic complexity = # of operations



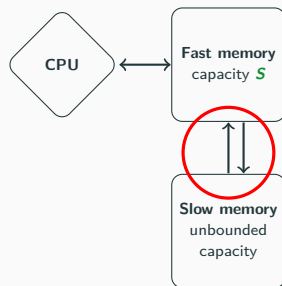
What is I/O complexity?

- Arithmetic complexity = # of operations



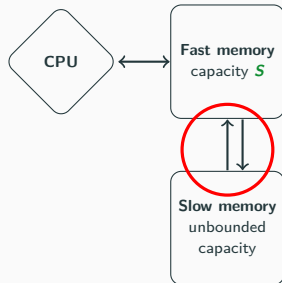
What is I/O complexity?

- Arithmetic complexity = # of operations
- I/O cost (schedule-dependent) = amount of data moved between fast and slow memory



What is I/O complexity?

- Arithmetic complexity = # of operations
- I/O cost (schedule-dependent) = amount of data moved between fast and slow memory
- I/O complexity = minimum cost over all schedules



Lower and Upper Bounds



IOLB (PLDI '20)
Automated lower
bound computation

Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs*

Jagan Srinivasan
University of Illinois at Urbana-Champaign
Urbana, IL, USA
jsriniva@illinois.edu

John Lapinskas
University of Illinois at Urbana-Champaign
Urbana, IL, USA
lapinska@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Pratik Panchbhaisankar
University of Illinois at Urbana-Champaign
Urbana, IL, USA
panchb1@illinois.edu

Lower and Upper Bounds



IOLB (PLDI '20)
Automated lower
bound computation

IOOpt (This paper)

- Improvement of the lower bound algorithm
- Automated upper bound derivation (IOUB)

Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs*

Joseph Chay
University of Illinois at Chicago
Joseph.Chay@uic.edu

John Lapinskas
University of Illinois at Chicago
John.Lapinskas@uic.edu

Levent N. Gurbuzoglu
University of Illinois at Chicago
Levent.Gurbuzoglu@uic.edu

Pratik P. Kamath
University of Illinois at Chicago
Pratik.Kamath@uic.edu

Shantanu Das
University of Illinois at Chicago
Shantanu.Das@uic.edu

Abstract
Memory movement is a critical component of program performance. In this paper, we present a novel algorithm for automated derivation of parametric lower bounds on the number of memory accesses required to execute an affine program. Our algorithm is based on a novel technique for deriving lower bounds on the number of memory accesses required to execute an affine program. Our algorithm is based on a novel technique for deriving lower bounds on the number of memory accesses required to execute an affine program.

1. Introduction
Memory movement is a critical component of program performance. In this paper, we present a novel algorithm for automated derivation of parametric lower bounds on the number of memory accesses required to execute an affine program. Our algorithm is based on a novel technique for deriving lower bounds on the number of memory accesses required to execute an affine program.

2. Preliminaries
In this section, we define the notation used in the paper. We use \mathbb{Z} to denote the set of integers, \mathbb{R} to denote the set of real numbers, and \mathbb{C} to denote the set of complex numbers. We use \mathbb{N} to denote the set of natural numbers, and \mathbb{Q} to denote the set of rational numbers. We use \mathbb{Z}^+ to denote the set of positive integers, and \mathbb{R}^+ to denote the set of positive real numbers. We use \mathbb{C}^+ to denote the set of positive complex numbers. We use $\mathbb{Z}^{\geq 0}$ to denote the set of non-negative integers, and $\mathbb{R}^{\geq 0}$ to denote the set of non-negative real numbers. We use $\mathbb{C}^{\geq 0}$ to denote the set of non-negative complex numbers. We use $\mathbb{Z}^{\leq 0}$ to denote the set of non-positive integers, and $\mathbb{R}^{\leq 0}$ to denote the set of non-positive real numbers. We use $\mathbb{C}^{\leq 0}$ to denote the set of non-positive complex numbers. We use $\mathbb{Z}^{\neq 0}$ to denote the set of non-zero integers, and $\mathbb{R}^{\neq 0}$ to denote the set of non-zero real numbers. We use $\mathbb{C}^{\neq 0}$ to denote the set of non-zero complex numbers. We use $\mathbb{Z}^{\neq 0}$ to denote the set of non-zero integers, and $\mathbb{R}^{\neq 0}$ to denote the set of non-zero real numbers. We use $\mathbb{C}^{\neq 0}$ to denote the set of non-zero complex numbers.

3. Problem Statement
In this section, we define the problem statement. We consider an affine program P with n memory locations. The program P is represented by a set of memory accesses \mathcal{M} . Each memory access $m \in \mathcal{M}$ is represented by a vector $\mathbf{m} \in \mathbb{Z}^n$. The program P is represented by a set of memory accesses \mathcal{M} . Each memory access $m \in \mathcal{M}$ is represented by a vector $\mathbf{m} \in \mathbb{Z}^n$. The program P is represented by a set of memory accesses \mathcal{M} . Each memory access $m \in \mathcal{M}$ is represented by a vector $\mathbf{m} \in \mathbb{Z}^n$.

4. Algorithm
In this section, we describe the algorithm. The algorithm is based on a novel technique for deriving lower bounds on the number of memory accesses required to execute an affine program. Our algorithm is based on a novel technique for deriving lower bounds on the number of memory accesses required to execute an affine program.

5. Results
In this section, we present the results of our algorithm. We show that our algorithm is able to derive lower bounds on the number of memory accesses required to execute an affine program. Our algorithm is based on a novel technique for deriving lower bounds on the number of memory accesses required to execute an affine program.

6. Conclusion
In this section, we conclude the paper. We have presented a novel algorithm for automated derivation of parametric lower bounds on the number of memory accesses required to execute an affine program. Our algorithm is based on a novel technique for deriving lower bounds on the number of memory accesses required to execute an affine program.

7. Acknowledgments
We thank the anonymous reviewers for their helpful comments. This work was supported by the National Science Foundation (NSF) grant CCF-1510045.

8. References
[1] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

[2] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

[3] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

[4] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

[5] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

[6] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

[7] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

[8] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

[9] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

[10] J. Chay, J. Lapinskas, L. N. Gurbuzoglu, P. P. Kamath, and S. Das, "Automated Derivation of Parametric Data Movement Lower Bounds for Affine Programs," *PLDI '20*, pp. 1–11, 2020.

I/O complexity upper bound \Leftrightarrow Cost of a particular valid schedule

Untiled matrix multiplication

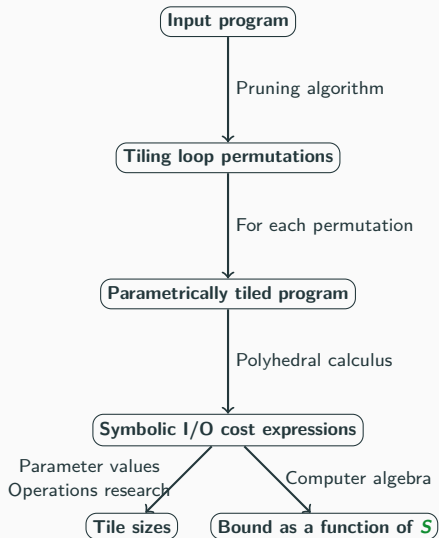
I/O cost: $O(N^3)$

Tiled matrix multiplication

I/O cost: $O(\frac{N^3}{\sqrt{5}})$

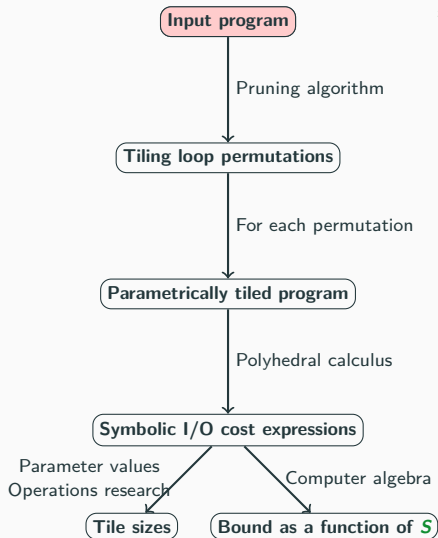
→ How to automatically compute I/O cost for a given schedule?

Upper bound derivation

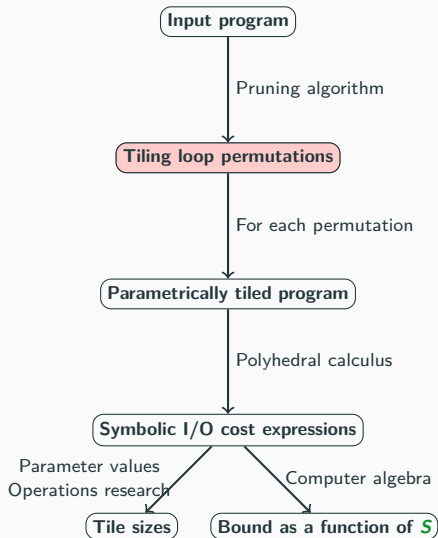


Upper bound derivation

```
for(i = 0; i < Ni; i++)  
  for(j = 0; j < Nj; j++)  
    for(k = 0; k < Nk; k++)  
      C[i][j] += A[i][k] * B[k][j];
```



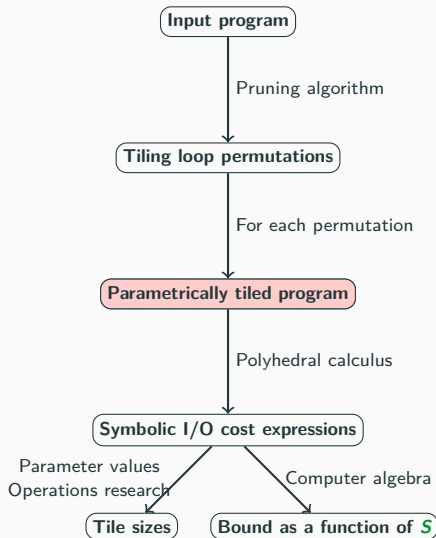
Upper bound derivation



```
for(i = 0; i < Ni; i++)
  for(j = 0; j < Nj; j++)
    for(k = 0; k < Nk; k++)
      C[i][j] += A[i][k] * B[k][j];
```

$\{(i, j, k), (i, k, j), (k, j, i)\}$

Upper bound derivation

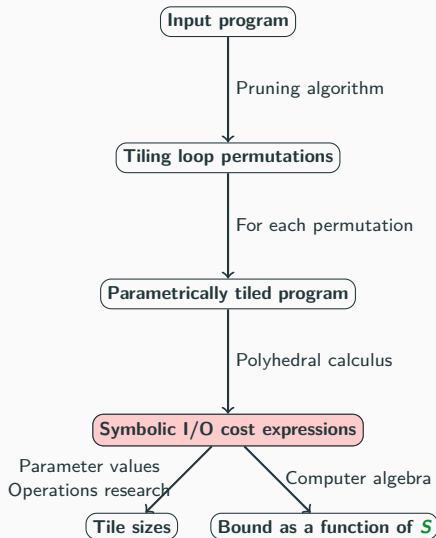


```
for(i = 0; i < Ni; i++)  
  for(j = 0; j < Nj; j++)  
    for(k = 0; k < Nk; k++)  
      C[i][j] += A[i][k] * B[k][j];
```

$\{(i, j, k), (i, k, j), (k, j, i)\}$

```
for(i1 = 0; i1 < Ni; i1+=Ti)  
  for(j1 = 0; j1 < Nj; j1+=Tj)  
    for(k = 0; k < Nk; k++)  
      for(i = i1; i < i1+Ti; i++)  
        for(j = j1; j < j1+Tj; j++)  
          C[i][j] += A[i][k] * B[k][j];
```

Upper bound derivation



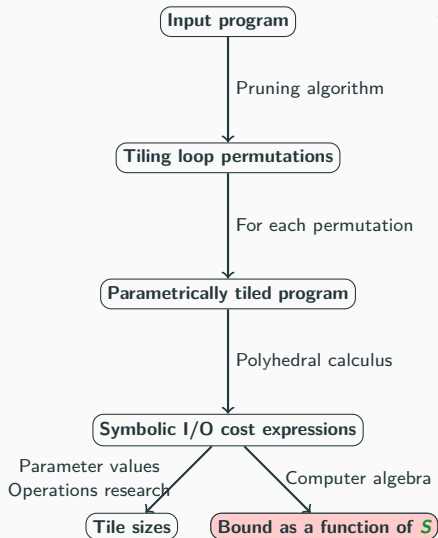
```
for(i = 0; i < Ni; i++)
  for(j = 0; j < Nj; j++)
    for(k = 0; k < Nk; k++)
      C[i][j] += A[i][k] * B[k][j];
```

$\{(i, j, k), (i, k, j), (k, j, i)\}$

```
for(i1 = 0; i1 < Ni; i1+=Ti)
  for(j1 = 0; j1 < Nj; j1+=Tj)
    for(k = 0; k < Nk; k++)
      for(i = i1; i < i1+Ti; i++)
        for(j = j1; j < j1+Tj; j++)
          C[i][j] += A[i][k] * B[k][j];
```

$$IO = N_i N_j N_k \left(\frac{1}{T_i} + \frac{1}{T_j} + \frac{1}{N_k} \right)$$
$$T_i T_j + T_i + T_j \leq S$$

Upper bound derivation



```
for(i = 0; i < Ni; i++)  
  for(j = 0; j < Nj; j++)  
    for(k = 0; k < Nk; k++)  
      C[i][j] += A[i][k] * B[k][j];
```

$\{(i, j, k), (i, k, j), (k, j, i)\}$

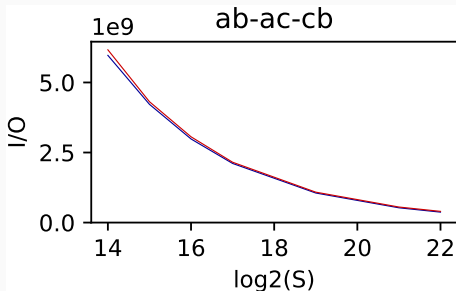
```
for(i1 = 0; i1 < Ni; i1+=Ti)  
  for(j1 = 0; j1 < Nj; j1+=Tj)  
    for(k = 0; k < Nk; k++)  
      for(i = i1; i < i1+Ti; i++)  
        for(j = j1; j < j1+Tj; j++)  
          C[i][j] += A[i][k] * B[k][j];
```

$$IO = N_i N_j N_k \left(\frac{1}{T_i} + \frac{1}{T_j} + \frac{1}{N_k} \right)$$
$$T_i T_j + T_i + T_j \leq S$$

$$UB = N_i N_j \left(\frac{2N_k}{\sqrt{S+1}-1} + 1 \right)$$

Matrix multiplication I/O complexity

$$\frac{2N_i N_j (N_k - 1)}{\sqrt{S}} \leq IO_{mm} \leq \frac{2N_i N_j N_k}{\sqrt{S+1}-1}$$



In the paper: Analytical results on several convolutions (Yolo9000) and tensor contractions (TCCG), with matching lower and upper bounds

Thank you!