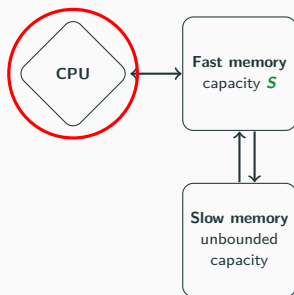


IOOpt: Automatic Derivation of I/O Complexity Bounds for Affine Programs

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

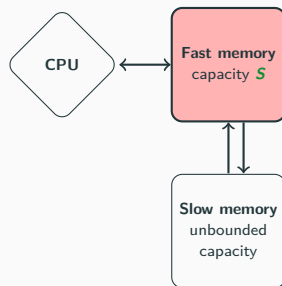
What is I/O complexity?

- Arithmetic complexity = # of operations



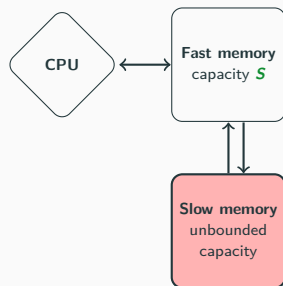
What is I/O complexity?

- Arithmetic complexity = # of operations



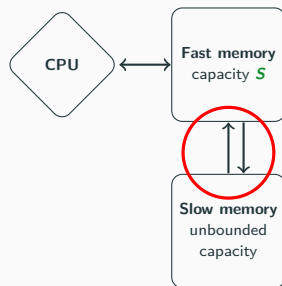
What is I/O complexity?

- Arithmetic complexity = # of operations



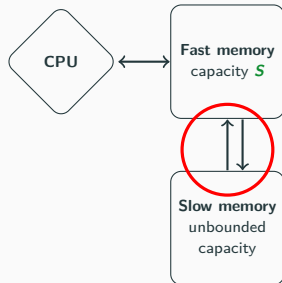
What is I/O complexity?

- 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*

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

Shihui Luo
University of Illinois at Urbana-Champaign
shihui@uiuc.edu

Amotz Oren
Technion, Haifa, Israel
amotz@technion.ac.il

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 Urbana-Champaign
Joseph.Chay@illinois.edu

John Lapinskas
University of Illinois at Urbana-Champaign
John.Lapinskas@illinois.edu

Levent N. Gurbuzoglu
University of Illinois at Urbana-Champaign
Levent.Gurbuzoglu@illinois.edu

Pratik P. Kamath
University of Illinois at Urbana-Champaign
Pratik.Kamath@illinois.edu

Chaitin S. Ramesh
University of Illinois at Urbana-Champaign
Chaitin.Ramesh@illinois.edu

Abstract
Memory and processor caches are key components of modern computer architectures. Understanding the memory and processor cache usage of a program is essential for understanding its performance. In this paper, we present a novel algorithm for automatically deriving parametric lower bounds on the number of memory accesses and cache misses of affine programs. Our algorithm is based on a novel technique for automatically deriving parametric lower bounds on the number of memory accesses and cache misses of affine programs. Our algorithm is based on a novel technique for automatically deriving parametric lower bounds on the number of memory accesses and cache misses of affine programs.

1. Introduction
Memory and processor caches are key components of modern computer architectures. Understanding the memory and processor cache usage of a program is essential for understanding its performance. In this paper, we present a novel algorithm for automatically deriving parametric lower bounds on the number of memory accesses and cache misses of affine programs. Our algorithm is based on a novel technique for automatically deriving parametric lower bounds on the number of memory accesses and cache misses of affine programs. Our algorithm is based on a novel technique for automatically deriving parametric lower bounds on the number of memory accesses and cache misses of affine programs.

2. Preliminaries
In this section, we review some basic concepts and notation used in this paper. We assume that the reader is familiar with the basics of affine programs and memory access patterns.

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

Untiled matrix multiplication

```
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];
```

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

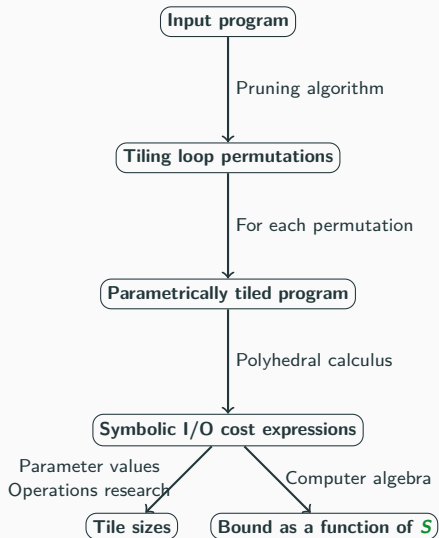
Tiled matrix multiplication

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

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

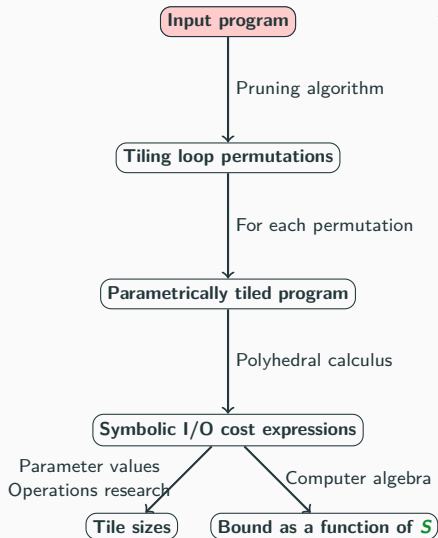
→ How to automatically compute I/O cost for a given schedule?

Upper bound derivation

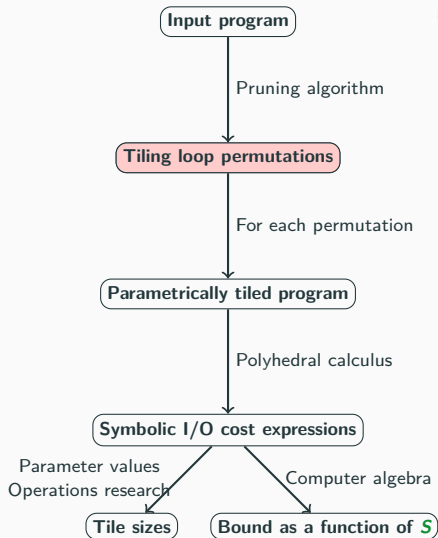


Upper bound derivation

```
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];
```



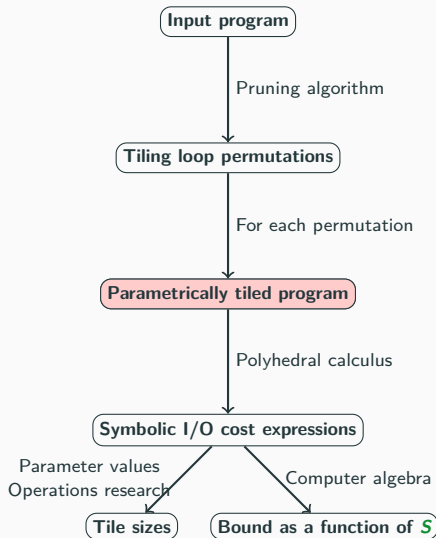
Upper bound derivation



```
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

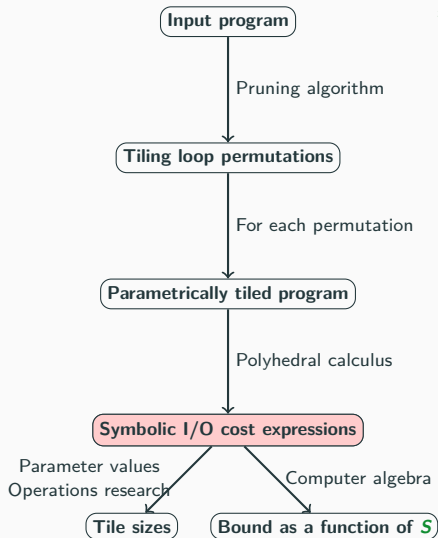


```
for(i = 0; i < Ni; i++)  
  for(j = 0; j < Nj; j++)  
    for(k = 0; k < Nk; k++)  
      C[i][j] += A[i][k] * B[k][j];
```

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

```
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



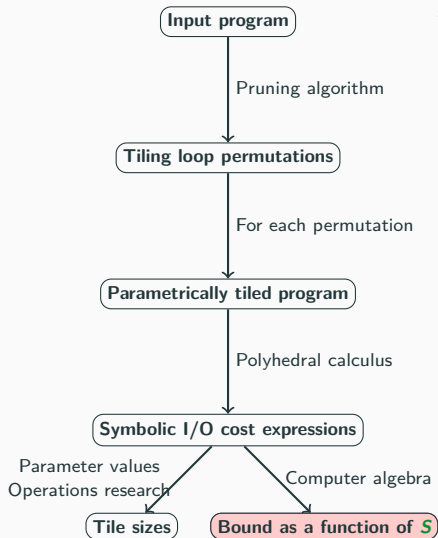
```
for(i = 0; i < Ni; i++)  
  for(j = 0; j < Nj; j++)  
    for(k = 0; k < Nk; k++)  
      C[i][j] += A[i][k] * B[k][j];
```

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

```
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];
```

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

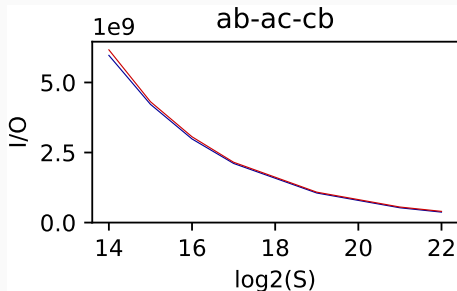
```
for(i1 = 0; i1 < Ni; i1+=Ti)
  for(j1 = 0; j1 < Nj; j1+=Tj)
    for(k = 0; k < Nk; k++)
      for(i = i1; i < i1+Ti; i++)
        for(j = j1; j < j1+Tj; j++)
          C[i][j] += A[i][k] * B[k][j];
```

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

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

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

TTile: Highly Optimized Tensor Computations

- Multi-level tiling driven by IOOpt model
- Microkernel: highly tuned “basic block” (vectorization, register reuse, instruction-level parallelism)

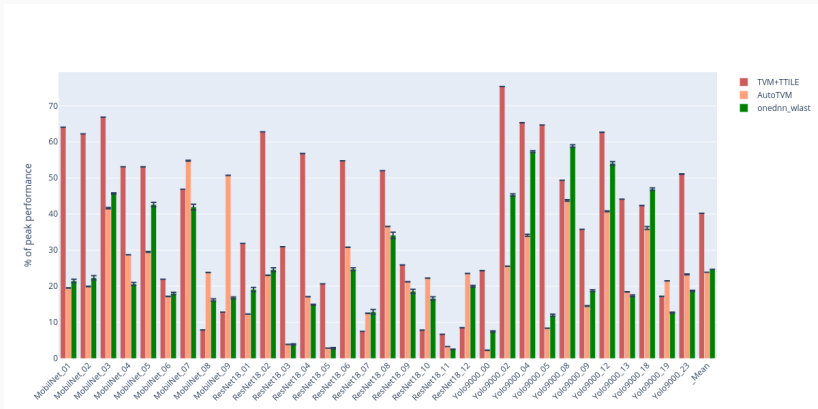
I/O: model-driven tiling (IOOpt)

```
for (kt = 0; kt < Nk; kt+=128)  
  for (it = 0; it < Ni; it+=32)  
    for (jt = 0; jt < Nj; jt+=6)
```

```
      for (k = kt; k < kt+128; k++)  
         $\mu$ kernel_gemm(A, B, C, i1, j1, k)
```

CPU: microkernel selection

TTile: Highly Optimized Tensor Contractions



Performance comparison between AutoTVM, oneDNN, and TTile+TVM for AVX512 (Intel Xeon Gold 6130), shown as percentage of machine peak. 32 threads were used, no hyperthreading

Demo

Thank you!