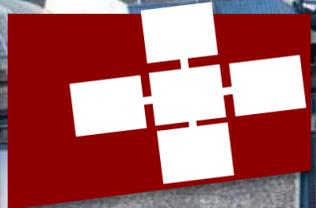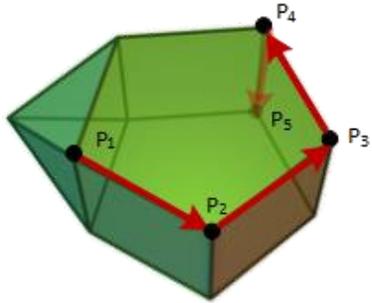# Fast Cache Modeling using Polyhedral Compilation

Tobias Gysi, **Tobias Grosser,** Laurin Brandner, Torsten Hoefler

# Historic Overview of Polyhedral Compilation

Mathematical
Foundations

Basic
Research

```
    int sum = 0;
    for(int i=0; i<4; ++i)
S0:   M[i] = i;
    for(int j=0; j<4; ++j)
S1:   sum += M[3-j];
```

Production Compilers
GCC/Graphite,
LLVM/Polly
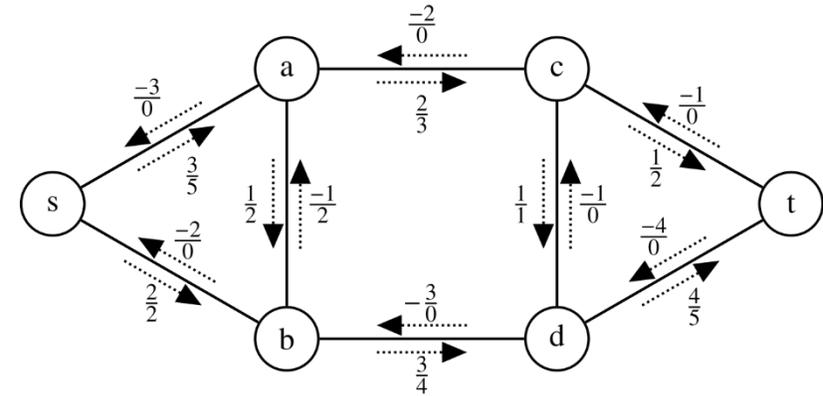
1960 - 1990

1990 - 2000

2000 - 2015

2015 - 2019

# Polyhedral Compilation: Analyses to build your own Magic



Memory
Footprints

Data-Transfer
Volumes

Data Dependences

# The Cost of Data Movement Depend on Global State and Do Not Compose

```
int N = 1000;
for(int i = 0; i < N; i++) {
  for(int j = 0; j < i; j++) {
    for(int k = 0; k < j; k++) {
      A[i][j] -= A[i][k] * A[j][k];
    }
    A[i][j] /= A[j][j];
  }
  for(int k = 0; k < i; k++) {
    A[i][i] -= A[i][k] * A[i][k];
  }
  A[i][i] = sqrt(A[i][i]);
}
```

Cholesky kernel from https://sourceforge.net/projects/polybench/

- cache sizes (32k and 512k)
- cacheline size (64B)

percentage of cache misses?

L1 cache **1.6%**

L2 cache **1.4%**

most expensive memory access?

**A[j][k]**

amount of compulsory and capacity misses?

| | |
|---|---|
| # compulsory misses | **31,752** |
| # capacity misses | **10,630,620** |

# HayStack Output for Cholesky Factorization

relative number of cache misses (statement)

parameters:
- cache sizes (32k and 512k)
- cacheline size (64B)

```
5    for (int i = 0; i < N; i++) {
6      for (int j = 0; j < i; j++) {
7        for (int k = 0; k < j; k++) {
8          A[i][j] -= A[i][k] * A[j][k];
------------------------------------------------------------------------
           ref    type  comp[%]     L1[%]      L2[%]     tot[%]     reuse[ln]
         A[i][j]  rd    0.00459    0.00000    0.00000   24.86910   8,10
         A[i][k]  rd    0.00000    0.00000    0.00000   24.86910   8,10
         A[j][k]  rd    0.00000    1.58635    1.38213   24.86910   8,10,13,15
         A[i][j]  wr    0.00000    0.00000    0.00000   24.86910   8
------------------------------------------------------------------------
```
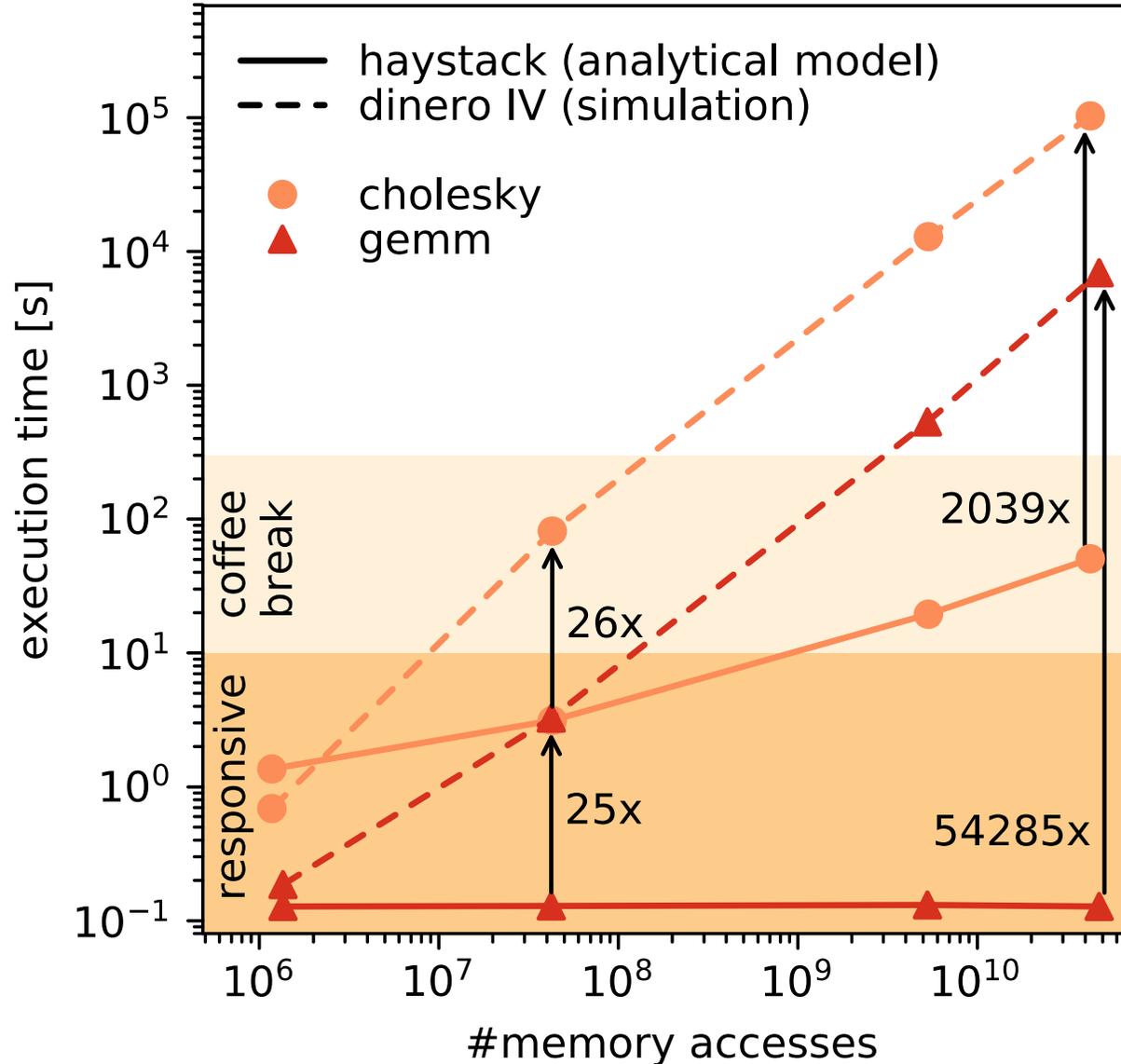
absolute number of cache misses (program)

```
compulsory:                       31'752
capacity (L1):                10'630'620
capacity (L2):                 9'258'460
total:                       668'166'500
```
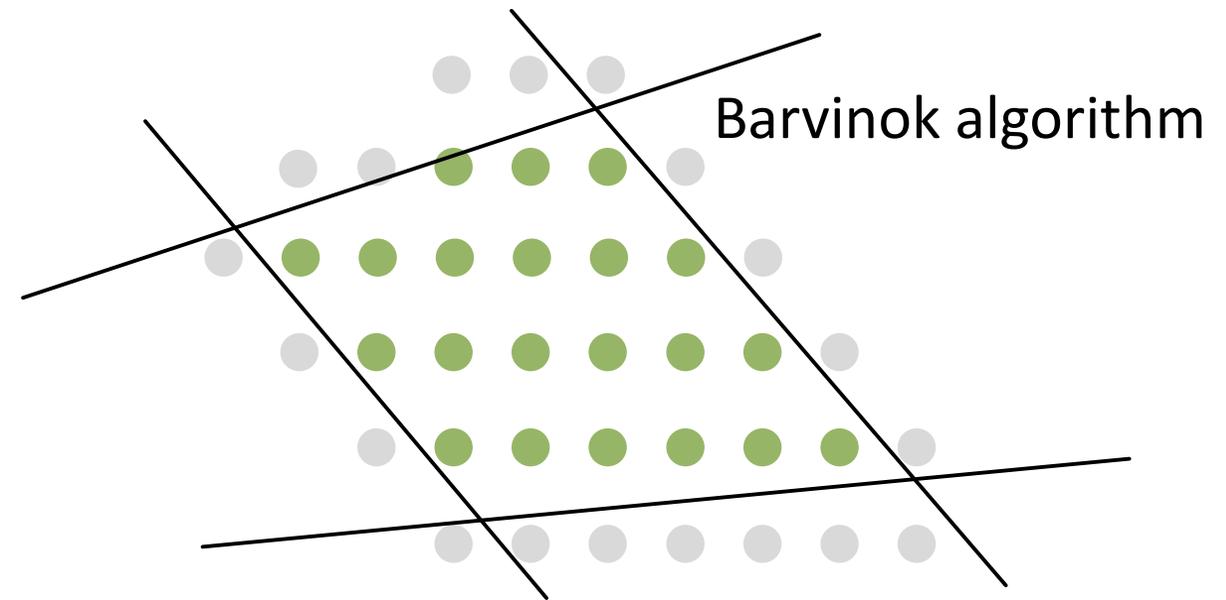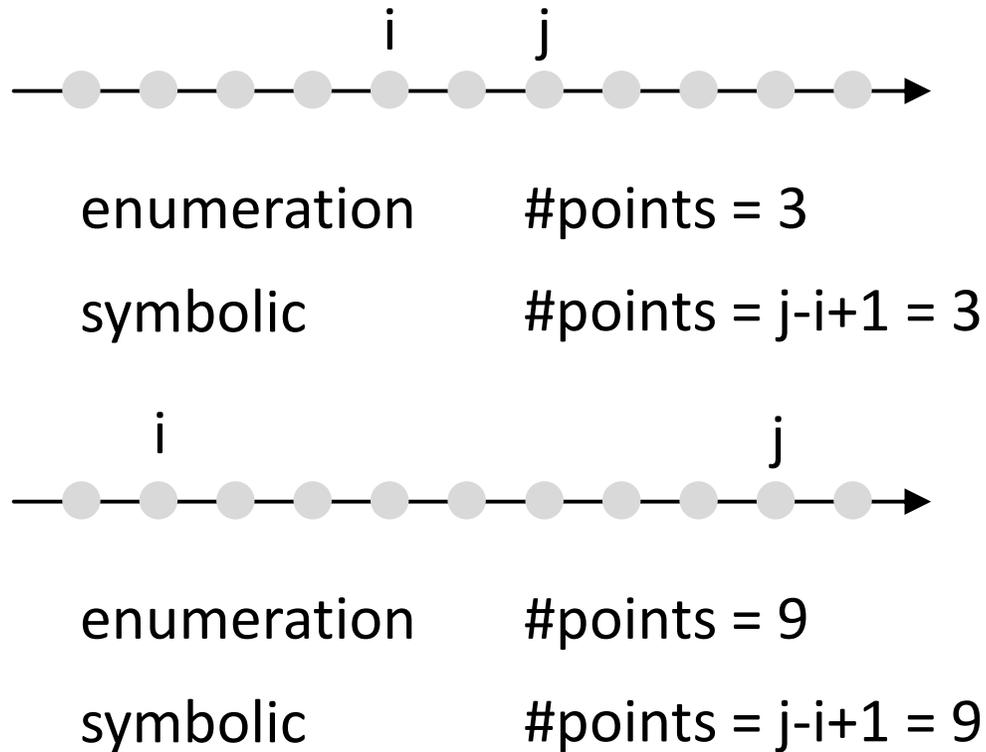
# Comparison to Simulation

# Symbolic Counting Avoids the Explicit Enumeration

1d illustration

i        j

enumeration        #points = 3

symbolic        #points = j-i+1 = 3

i                j

Barvinok algorithm

enumeration        #points = 9

symbolic        #points = j-i+1 = 9

Alexander I. Barvinok, *A Polynomial Time Algorithm for Counting Integral Points in Polyhedra When the Dimension is Fixed*. 1994.

# The LRU Stack Distance Allows Us to Model Fully Associative Caches

example

```
    int sum = 0;
    for(int i=0; i<4; ++i)
S0:    M[i] = i;
    for(int j=0; j<4; ++j)
S1:    sum += M[3-j];
```
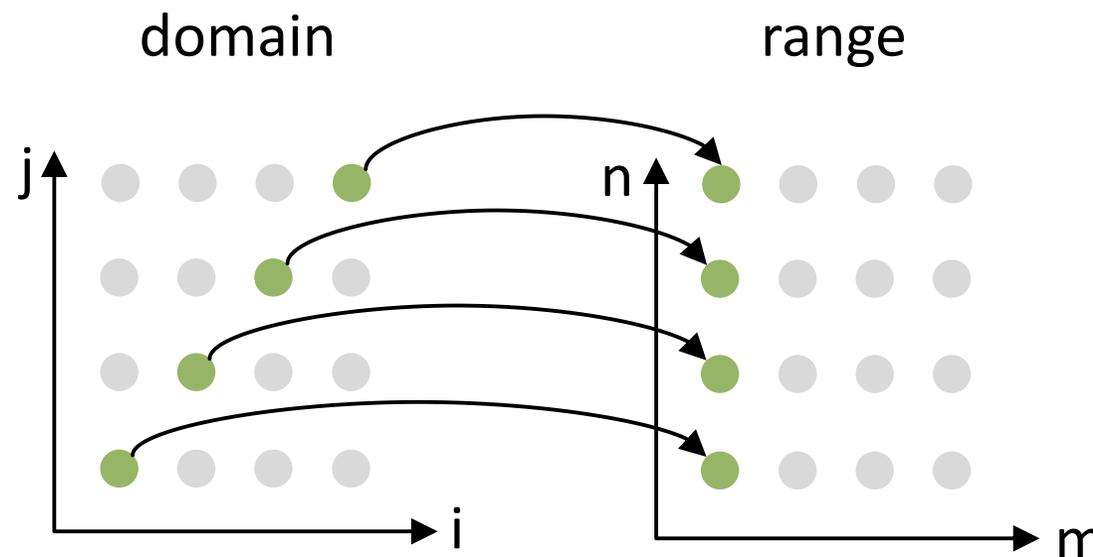
deliberately **generic model**

Richard L Mattson, Jan Gecsei, Donald R Slutz, and Irving L Traiger, *Evaluation techniques for storage hierarchies.* 1970.

# Integer Sets and Relations
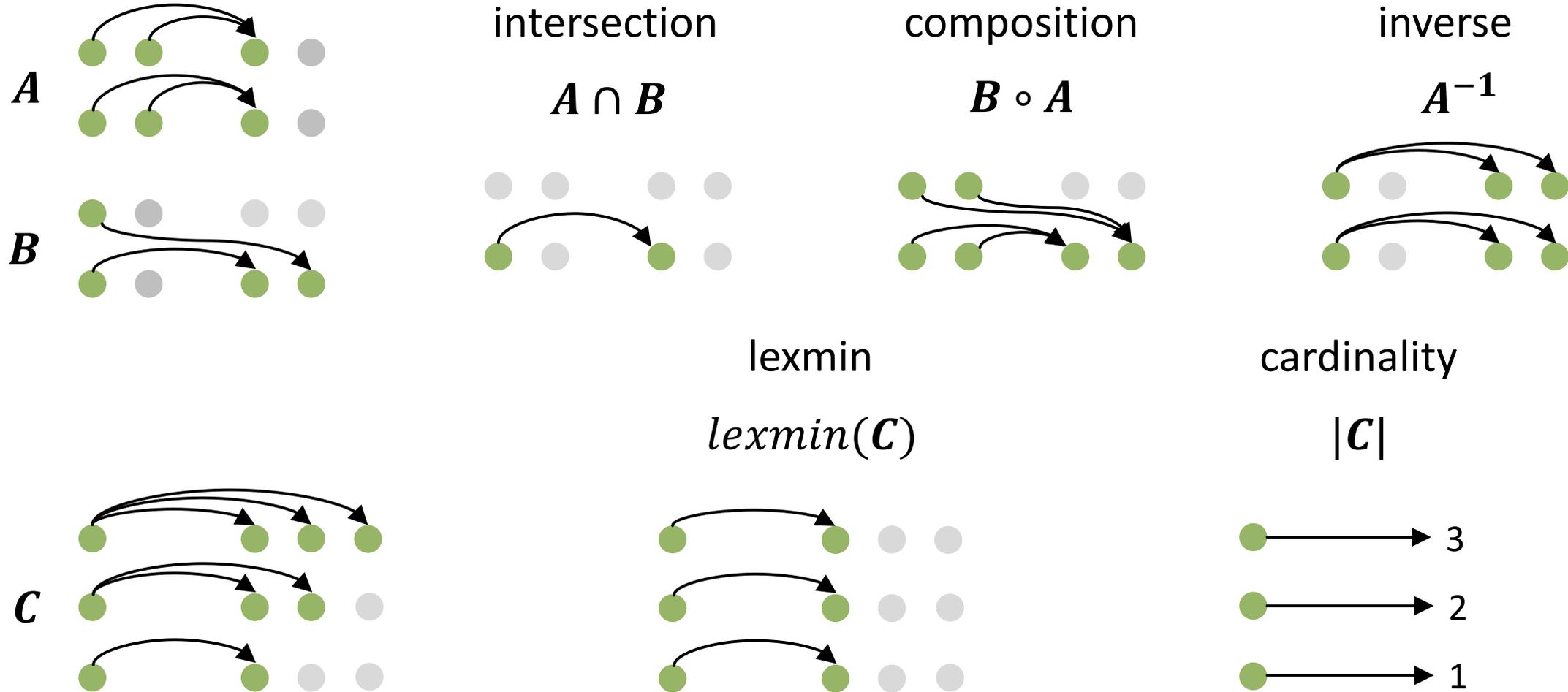
affine set

affine relation

domain                    range



$$\mathbf{S} = \{(i,j) : 0 < i < 3 \wedge 0 < j < 3\} \quad \mathbf{R} = \{(i,j) \rightarrow (m,n) : 0 \leq i < 4 \wedge i = j \wedge n = j \wedge m = 0\}$$

Sven Verdoolaege, *isl: An integer set library for the polyhedral model.* 2010.

# Computation with Affine Sets and Relations



intersection

$$A \cap B$$

composition

$$B \circ A$$

inverse

$$A^{-1}$$

lexmin

$$lexmin(C)$$

cardinality

$$|C|$$

Sven Verdoolaege, *isl: An integer set library for the polyhedral model.* 2010.

# Express Programs by Integer Sets and Relations

example

schedule values    statement instances    memory accesses

```
        int sum = 0;
        for(int i=0; i<4; ++i)
S0:     M[i] = i;

        for(int j=0; j<4; ++j)
S1:     sum += M[3-j];
```

i=0    (0,0) ←———————— S0(0) ————————→ M(0)

i=1    (0,1) ←———————— S0(1) ————————→ M(1)

i=2    (0,2) ←———————— S0(2) ————————→ M(2)

i=3    (0,3) ←———————— S0(3) ————————→ M(3)

j=0    (1,0) ←———————— S1(0)

j=1    (1,1) ←———————— S1(1)

j=2    (1,2) ←———————— S1(2)

j=3    (1,3) ←———————— S1(3)

*S*                                  *A*

# Express Programs by Integer Sets and Relations

schedule values

$$\{(i0, i1) : 0 \leq i0 < 2 \wedge$$
$$0 \leq i1 < 4\}$$

statement instances

$$\{S0(i) : 0 \leq i < 4;$$
$$S1(j) : 0 \leq j < 4\}$$

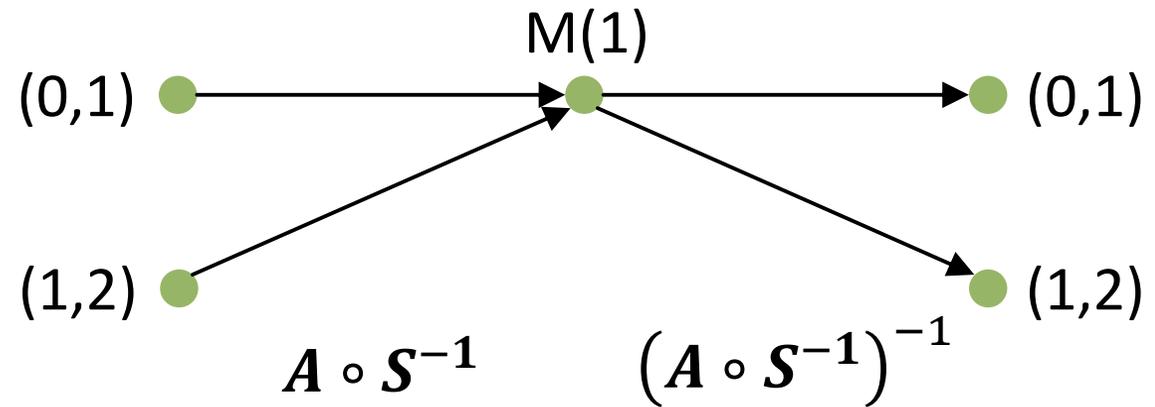memory accesses

$$\{M(k) : 0 \leq k < 4\}$$



$$\boldsymbol{S} = \{S0(i) \rightarrow (0, i) : 0 \leq i < 4;$$
$$S1(j) \rightarrow (1, j) : 0 \leq j < 4\}$$

$$\boldsymbol{A} = \{S0(i) \rightarrow M(k) : k = i;$$
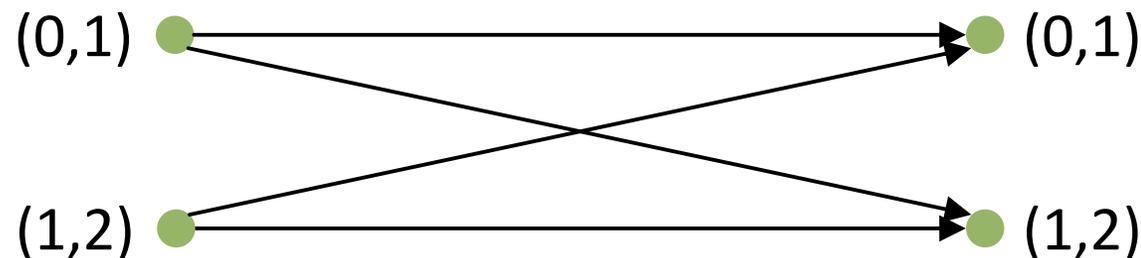$$S1(j) \rightarrow M(k) : k = 3 - j\}$$

# Identify Next Access of the Same Memory Location

example

```
int sum = 0;
for(int i=0; i<4; ++i)
S0:   M[i] = i;
for(int j=0; j<4; ++j)
S1:   sum += M[3-j];
```



M(1)

(0,1) → M(1) → (0,1)

(1,2) → M(1) → (1,2)

$$A \circ S^{-1} \qquad (A \circ S^{-1})^{-1}$$

# Identify Next Access of the Same Memory Location

example

```
    int sum = 0;
    for(int i=0; i<4; ++i)
S0:   M[i] = i;
    for(int j=0; j<4; ++j)
S1:   sum += M[3-j];
```

$$E = (A \circ S^{-1})^{-1} \circ A \circ S^{-1}$$

$$N = lexmin(E \cap L_{\prec})$$

(0,1) ●  ————————————→  ● (0,1)

(1,2) ●  ————————————→  ● (1,2)

$$lexmin(E \cap L_{\prec})$$

$$L_{\prec} = \{(i0, i1) \rightarrow (j0, j1) :$$
$$(i0 < j0) \vee$$
$$(i0 = j0 \wedge i1 < j1)\}$$

# Compute the LRU Stack Distance

example

```
    int sum = 0;
    for(int i=0; i<4; ++i)
S0:   M[i] = i;
    for(int j=0; j<4; ++j)
S1:   sum += M[3-j];
```
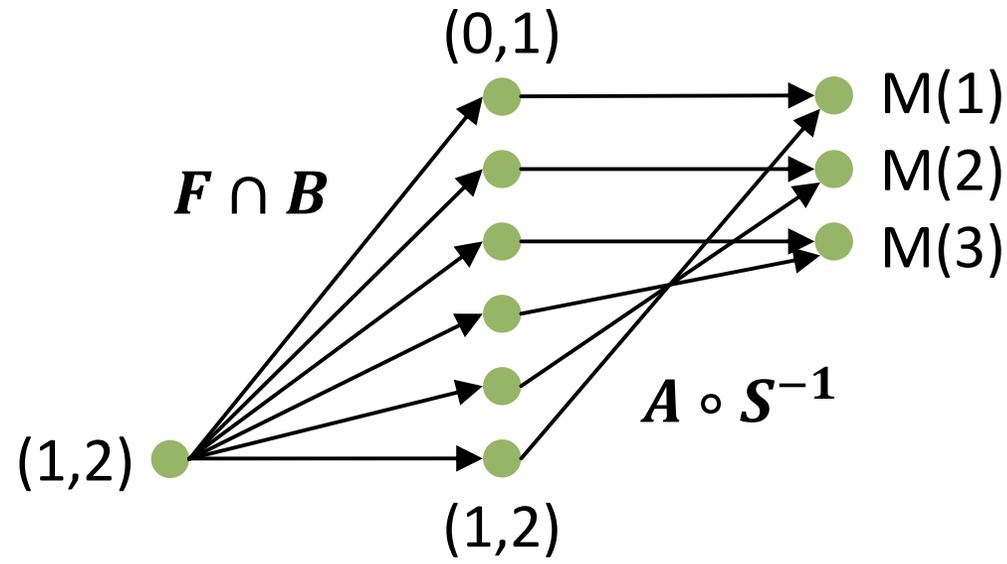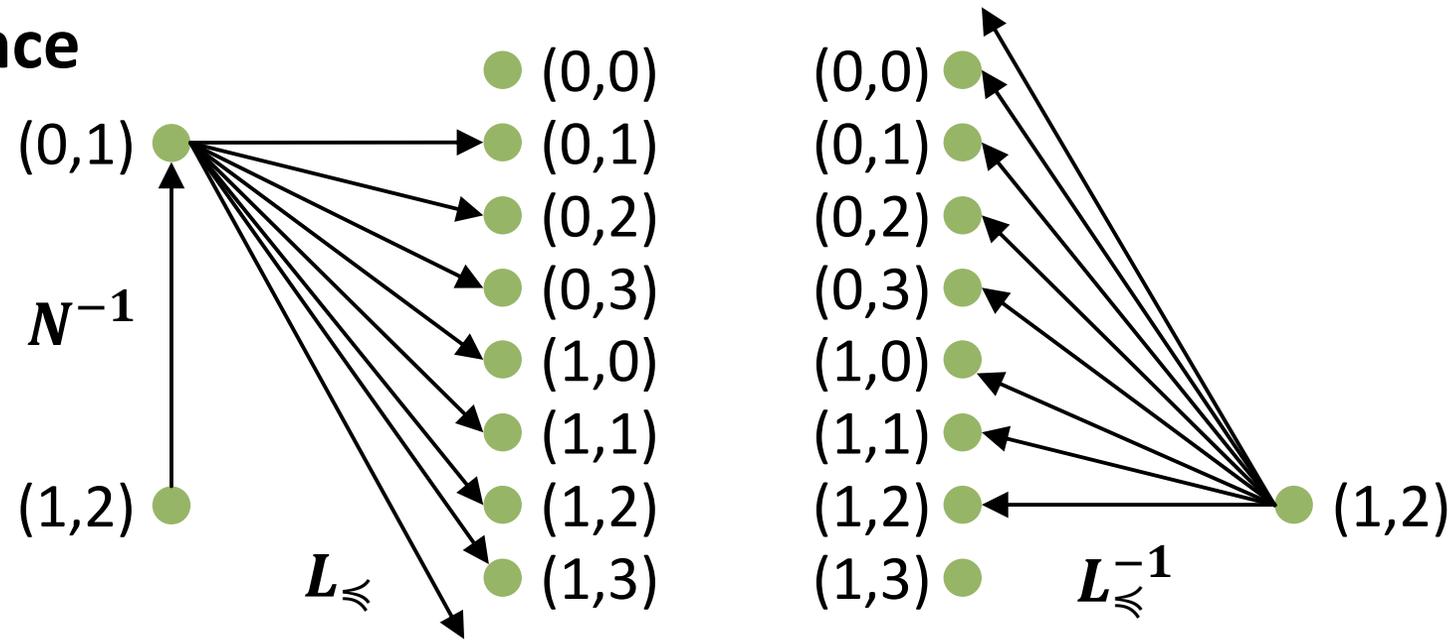
$$E = (A \circ S^{-1})^{-1} \circ A \circ S^{-1}$$

$$N = lexmin(E \cap L_{\prec})$$

$$F = L_{\preccurlyeq} \circ N^{-1}$$

$$B = L_{\preccurlyeq}^{-1}$$

$$D = \{|A \circ S^{-1} \circ (F \cap B) \circ S|\}$$

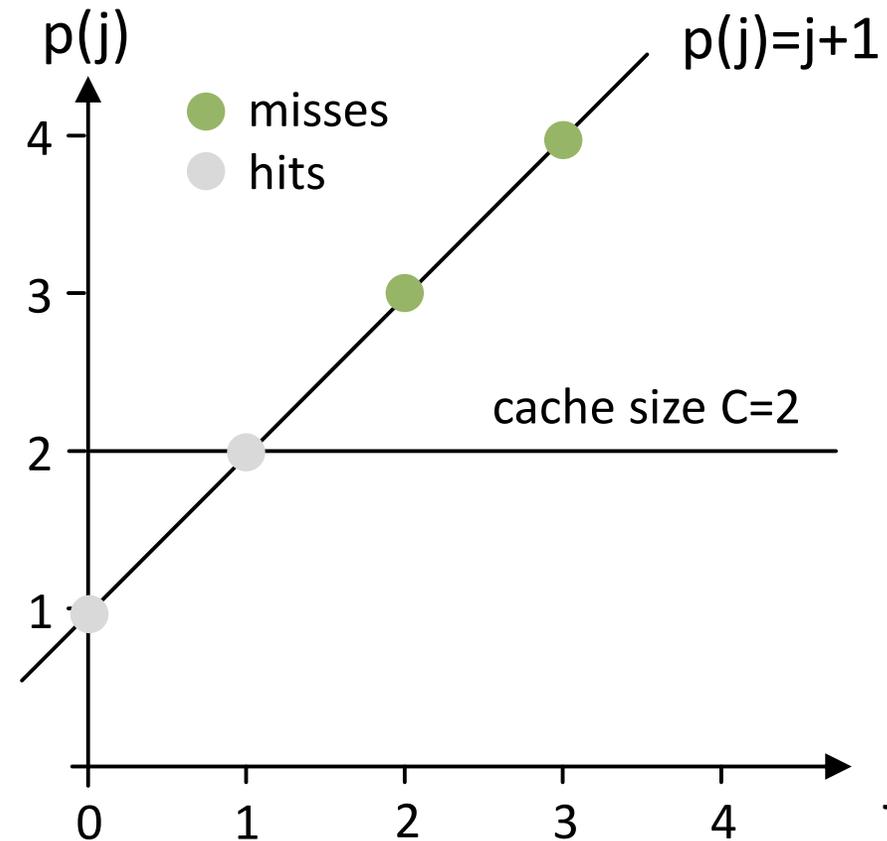# Count the Cache Misses Given the Backward LRU Stack Distance

example

```
int sum = 0;
for(int i=0; i<4; ++i)
S0:    M[i] = i;
for(int j=0; j<4; ++j)
S1:    sum += M[3-j];
```
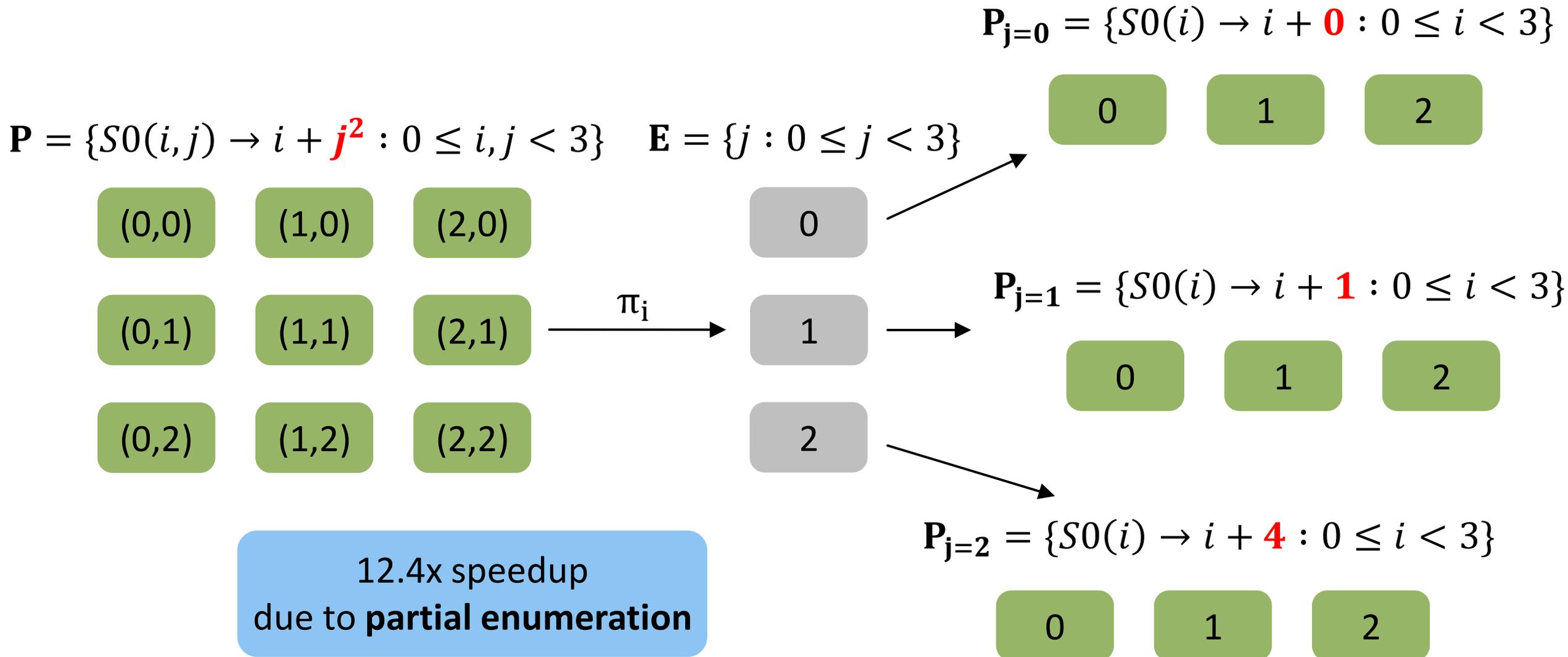
$$\mathbf{D} = \{|A \circ S^{-1} \circ (F \cap B) \circ S|\}$$

$$= \{S1(j) \to j + 1 : 0 \leq j < 4\}$$



$$|\mathbf{M}| = |\{S1(j) : p(j) > C \wedge 0 \leq j < 4\}|$$

$$= |\{S1(2), S1(3)\}| = 2$$

Kristof Beyls and Erik H D'Hollander, *Generating cache hints for improved program efficiency*. 2005.

16

# Enumerate the Non-Affine Dimensions of the Stack Distance Polynomials

$$\mathbf{P_{j=0}} = \{S0(i) \rightarrow i + \mathbf{\color{red}0} : 0 \leq i < 3\}$$

| 0 | 1 | 2 |
|---|---|---|

$$\mathbf{P} = \{S0(i,j) \rightarrow i + \mathbf{\color{red}j^2} : 0 \leq i, j < 3\} \quad \mathbf{E} = \{j : 0 \leq j < 3\}$$

| (0,0) | (1,0) | (2,0) |

| 0 |

| (0,1) | (1,1) | (2,1) | $\xrightarrow{\pi_i}$ | 1 | $\rightarrow$

$$\mathbf{P_{j=1}} = \{S0(i) \rightarrow i + \mathbf{\color{red}1} : 0 \leq i < 3\}$$

| 0 | 1 | 2 |
|---|---|---|

| (0,2) | (1,2) | (2,2) |

| 2 |

$$\mathbf{P_{j=2}} = \{S0(i) \rightarrow i + \mathbf{\color{red}4} : 0 \leq i < 3\}$$

12.4x speedup
due to **partial enumeration**

| 0 | 1 | 2 |
|---|---|---|

# Accuracy of HayStack for the L1 Cache of Our Test System



$\overline{e}^{(g)} = 0.6\%$

# Error of HayStack Compared to Simulation (Dinero IV)

HayStack
(fully associative)



HayStack chart annotation: ← due to associativity, $\overline{e}^{(g)} = 0.6\%$

Dinero IV
(fully associative)



Dinero IV chart annotation: ← due to associativity, $\overline{e}^{(g)} = 0.6\%$

Dinero IV
(8-way associative)



Dinero IV chart annotation: $\overline{e}^{(g)} = 0.4\%$

# Performance of HayStack for the Large Problem Size of PolyBench

# Performance of HayStack Compared to PolyCache and Dinero

Dinero IV

- simulator
- setup to simulate full associativity
- problem size dependent performance

370x speedup
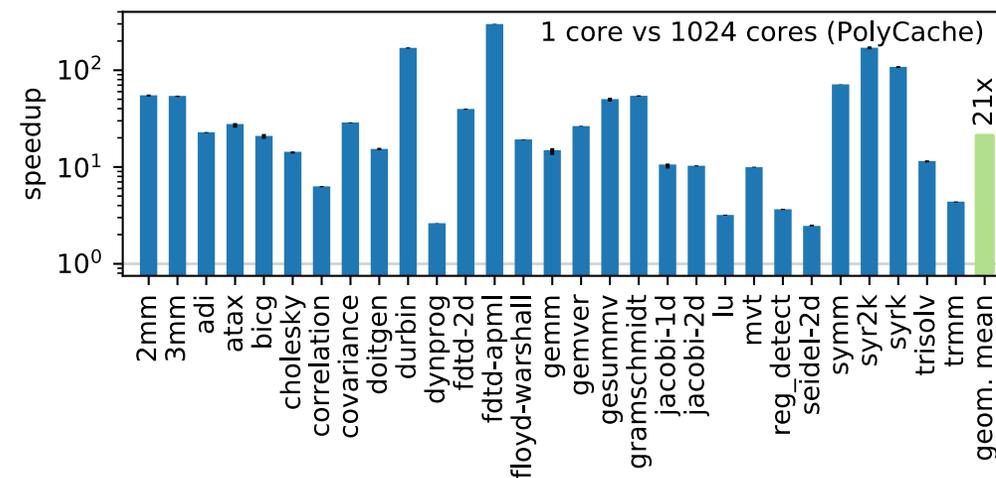


Jan Elder and Mark D. Hill, *Dinero IV Trace-Driven Uniprocessor Cache Simulator*. 2003.

PolyCache

- analytical cache model
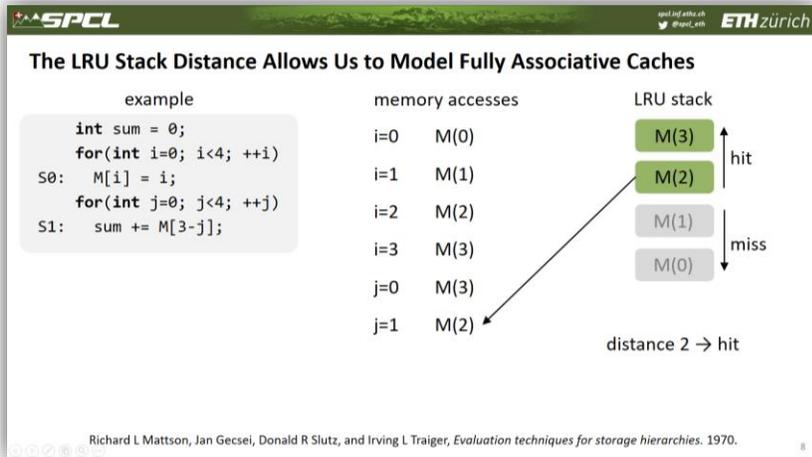- models set associativity
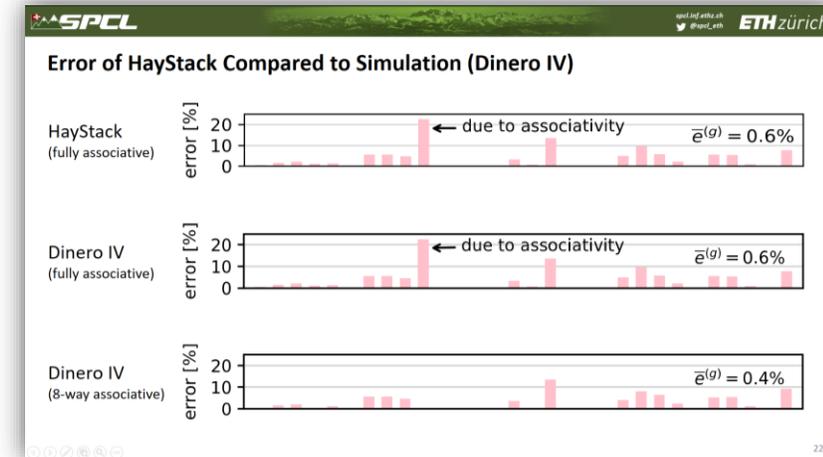- one core per cache set

21x speedup



Wenlei Bao, Sriram Krishnamoorthy, Louis-Noel Pouchet, and P Sadayappan, *Analytical modeling of cache behavior for affine programs.* 2017.
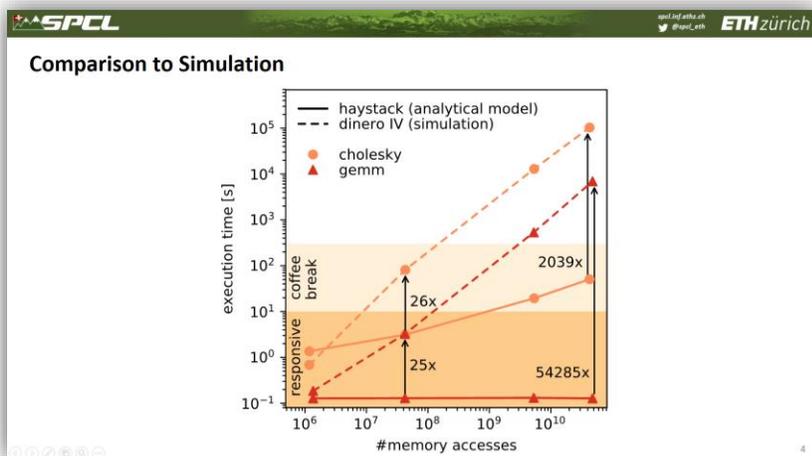
# Conclusion

### generic model of **fully associative caches**



### **accurate results** compared to measurements



### fast enough to provide **interactive feedback**



### **symbolic counting** avoids explicit enumeration