



Symbolic Footprint Analysis and Uses

Chen Ding

Professor

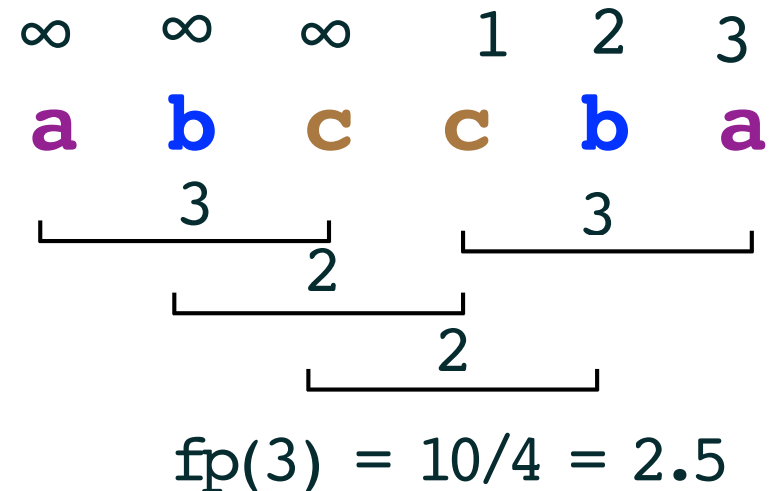
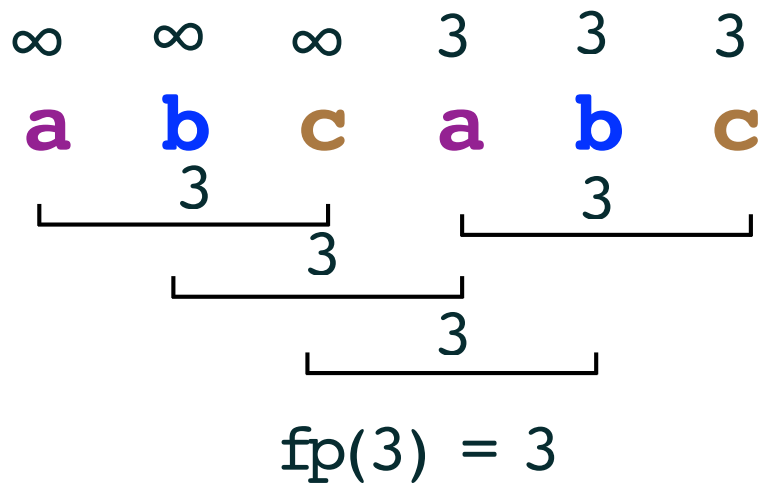
Department of Computer Science

University of Rochester

Locality Analysis

- What is locality?
 - Denning's definition
 - working set
 - how to quantify?
- Essence of cache
 - dynamic sharing of (fast) memory
 - memory allocation for working sets
 - within/across threads and programs
 - program optimization, programs management
- Locality analysis
 - program parameters
 - conversion to miss ratio (cache parameters)

- Locality of an access
 - shorter reuse distance \rightarrow better locality
- Locality of an execution window
 - smaller working set \rightarrow better locality
- Timescale locality: footprint $fp(x)$
 - average working set size of all windows of length x
 - parameterized by timescale
 - smaller footprint \rightarrow smaller average WSS \rightarrow better locality



Reuse Distance (rd) vs Reuse Time (rt)

- Distribution function $P(\text{rd})$ and $P(\text{rt})$

- e.g.

- $P(\text{rd}=3) = 2/7$

- $P(\text{rt}=3) = 0/7$

| | | | | | | | |
|-------------|----------|----------|----------|---|---|---|---|
| rt : | ∞ | ∞ | ∞ | 4 | 2 | 4 | 6 |
| rd : | ∞ | ∞ | ∞ | 3 | 1 | 2 | 3 |
| | a | b | c | a | a | c | b |

- $P(\text{rt})$ is easy to measure

- $P(\text{rt})$ can be used to compute footprint

- precise solution [Xiang et al. PACT'11]

- approximate solution, Denning 1968

- The higher-order theory of locality (HOTL)

- miss ratio, reuse distance, fill time [Xiang et al. ASPLOS'13]

- Growing literature [OSDI'14, FAST'15, ATC'16, MEMSYS'16]

Static Analysis of Footprint

- Building on reuse distance equations
 - Beyls and D'Hollander JSA 2004
 - nested loops and affine array subscript expressions
- Key concepts
 - references and iterations
 - reuse pairs
 - forward reuses
 - iteration polytopes

Step 1. Reuse Pairs

$$\forall r, s \in \mathcal{R} : \text{reuse}(r \rightarrow s) = \{(I_r, J_s) \in \mathbb{Z}^n \mid \text{subject to conditions (2a)--(2d)}\}, \quad (2)$$

$$I_r \in \text{IS}(r) \wedge J_s \in \text{IS}(s) \quad (\text{iteration space}), \quad (2a)$$

$$I_r \triangleleft J_s \quad (\text{execution order}), \quad (2b)$$

$$r@I_r = s@J_s \quad (\text{same location}), \quad (2c)$$

$$\forall t \in \mathcal{R} : \neg(\exists K_t \in \text{IS}(t) : I_r \triangleleft K_t \triangleleft J_s \wedge t@K_t = r@I_r) \quad (\text{no intervening access}). \quad (2d)$$

Step 2. Iteration Polytopes

```

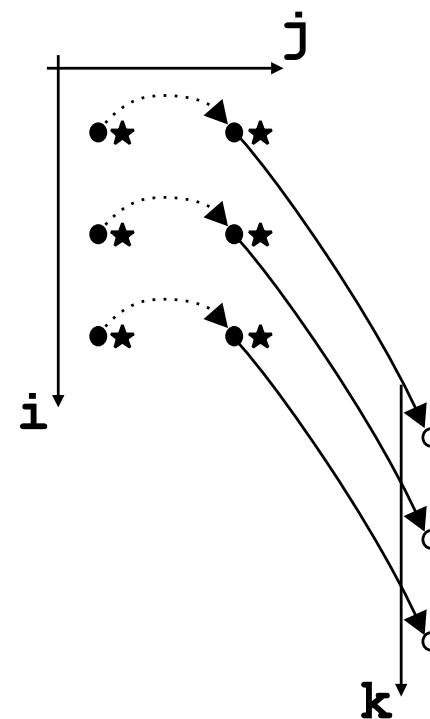
do i=1,N
  do j=1,2
    B(i,j) = A(i)
  enddo
enddo
do k=1,N
  A(k)=0
enddo

```

$$\text{reuse}(A(i) \rightarrow A(i)) = \{(i, j, i', j') \mid 1 \leq i, i' \leq N \wedge j = 1 \wedge j' = 2\}$$

$$\text{reuse}(A(i) \rightarrow A(k)) = \{(i, j, k) \mid 1 \leq i \leq N \wedge j = 2 \wedge k = i\}$$

$$\text{reuse}_F(A(i)) = \{(i, j) \mid 1 \leq i \leq N \wedge 1 \leq j \leq 2\}$$



..... reuse($A(i) \rightarrow A(i)$)

→ reuse($A(i) \rightarrow A(k)$)

$$\text{reuse}_F(r) = \bigcup_{\forall s \in \mathcal{R}} \{I_r \mid \exists J_s : (I_r, J_s) \in \text{reuse}(r \rightarrow s)\},$$

Static Reuse Distance



Bin Bao
[CGO 2013]

```

for(jj = 0; jj < N; jj = jj + B_j)
  for(kk = 0; kk < N; kk = kk + B_k)
    for(i = 0; i < N; i = i + 1)
      for(j = jj; j < min(jj + B_j, N); j = j + 1)
        for(k = kk; k < min(kk + B_k, N); k = k + 1)
          C[i][j] = beta * C[i][j] + alpha * A[i][k] * B[k][j];
  
```

| Loop | Reuse | Reuse distance (bytes) | Reuse time (accesses) |
|------|-----------|---|-----------------------|
| k | $C[i][j]$ | $8 * 3$ | 3 |
| j | $A[i][k]$ | $8 * 1 + 8 * B_j + 8 * B_j$ | $3 * B_j$ |
| i | $B[k][j]$ | $8 * B_k + 8 * B_k * B_j + 8 * B_j$ | $3 * B_j * B_k$ |
| jj | $A[i][k]$ | $8 * B_i * B_k + 8 * B_k * B_j + 8 * B_i * B_j$ | $3 * B_i * B_j * B_k$ |
| kk | $C[i][j]$ | $8 * B_i * B_k + 8 * B_k * J + 8 * B_i * J$ | $3 * B_i * J * B_k$ |
| ii | $B[k][j]$ | $8 * B_i * K + 8 * K * J + 8 * B_i * J$ | $3 * B_i * J * K$ |

Uses

- Locality optimization of GPU kernels
 - select the best order to execute the thread groups in a thread block
 - statically compute the footprint and miss ratio
 - static or dynamic optimization
 - select the best placement of arrays
 - statically compute the footprint and miss ratio
 - for every group of arrays
 - possible use in PORPLE [Chen et al. MICRO 2014]

Uses (cont'd)

- Optimization of NUMA task and data placement (TPDP)
 - joint work with Robert Ho and Yaoqing Gao
 - tomorrow's workshop
- Write locality
 - write caching [Li+ LCPC'16] in ATLAS [Chakrabarti+ OOPSLA'14]
 - write locality analysis / optimization [Chen+ MEMSYS'16]
 - modeling writebacks in mixed read-write cache
 - write reuse distance
 - fast measurement
 - compiler optimization of writes
- Other suggestions?