

Parametric Tiling with Inter-Tile Data Reuse

Alexandre Isoard Alain Darte

Comsys, LIP (Laboratoire de l'Informatique du Parallélisme), Lyon

IMPACT

4th International Workshop on
Polyhedral Compilation Techniques

January 20, 2014

Vienna, Austria



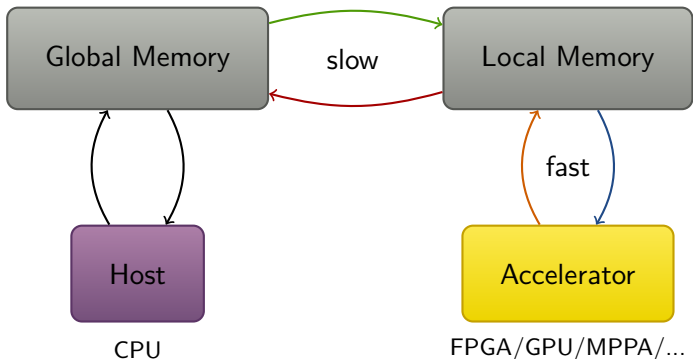
Lyon 1

UNIVERSITÉ DE LYON

Outline

- 1 Motivation and challenges
 - Kernel offloading: rules of the game
 - Reminders: scheduling and tiling
 - Inter-tile data reuse: example
- 2 Parametric analysis
 - Tile index vs tile origin index
 - Exact inter-tile reuse
 - Approximated inter-tile reuse
- 3 Current implementation and results
 - Current status
 - Script with ISCC
 - Local memory allocation for PolyBench examples

Kernel Offloading



- ☛ Perform computations by blocks;
- ☛ Exploit data reuse;
- ☛ Use pipelining/prefetching;
- ☛ Reduce and coalesce communications (burst).

Rules and objectives

Data reuse: on the full iteration domain

Rule 1: always use local data if already loaded or computed.

- Reduces communication volume, increases local memory.
- Enables full pipelining (load/compute/store sequence).

Rules and objectives

Data reuse: on the full iteration domain

Rule 1: always use local data if already loaded or computed.

- ☛ Reduces communication volume, increases local memory.
- ☛ Enables full pipelining (load/compute/store sequence).

Blocking: thanks to tiling

Rule 2: tiles executed in sequence (but a tile can be parallelized).

- ☛ Increases temporal reuse, reduces local memory.
- ☛ Increases spatial reuse, enables burst communications.

Rules and objectives

Data reuse: on the full iteration domain

Rule 1: always use local data if already loaded or computed.

- ☛ Reduces communication volume, increases local memory.
- ☛ Enables full pipelining (load/compute/store sequence).

Blocking: thanks to tiling

Rule 2: tiles executed in sequence (but a tile can be parallelized).

- ☛ Increases temporal reuse, reduces local memory.
- ☛ Increases spatial reuse, enables burst communications.

Variants for *reuse domain*, i.e., where data reuse is performed

- Iteration domain reduced thanks to hierarchical tiling.
- Data reuse in a p -dimensional stripe, or at bounded distance.

Rules and objectives

Data reuse: on the full iteration domain

Rule 1: always use local data if already loaded or computed.

- ☛ Reduces communication volume, increases local memory.
- ☛ Enables full pipelining (load/compute/store sequence).

Blocking: thanks to tiling

Rule 2: tiles executed in sequence (but a tile can be parallelized).

- ☛ Increases temporal reuse, reduces local memory.
- ☛ Increases spatial reuse, enables burst communications.

Variants for *reuse domain*, i.e., where data reuse is performed

- Iteration domain reduced thanks to hierarchical tiling.
- Data reuse in a p -dimensional stripe, or at bounded distance.

Then: scheduling/pipelining & memory allocation

Rule 3: reuse analysis independently on scheduling.

Rule 4: load as late as possible, store as soon as possible.

- ☛ Overlaps transfer and computation (multi-buffering).
- ☛ Reduces live-ranges, and possibly local memory size.

Rules and objectives

Parametric in terms of tile sizes?

Data reuse: on the full iteration domain

Rule 1: always use local data if already loaded or computed.

- ☛ Reduces communication volume, increases local memory.
- ☛ Enables full pipelining (load/compute/store sequence).

Blocking: thanks to tiling

Rule 2: tiles executed in sequence (but a tile can be parallelized).

- ☛ Increases temporal reuse, reduces local memory.
- ☛ Increases spatial reuse, enables burst communications.

Variants for *reuse domain*, i.e., where data reuse is performed

- Iteration domain reduced thanks to hierarchical tiling.
- Data reuse in a p -dimensional stripe, or at bounded distance.

Then: scheduling/pipelining & memory allocation

Rule 3: reuse analysis independently on scheduling.

Rule 4: load as late as possible, store as soon as possible.

- ☛ Overlaps transfer and computation (multi-buffering).
- ☛ Reduces live-ranges, and possibly local memory size.

Challenges and contributions

General principle for Load sets

Load a data indexed by \vec{m} just before a tile indexed by \vec{T} if:


- \vec{m} is live-in for \vec{T} , i.e., read but not written earlier in \vec{T} .
- \vec{m} has not been loaded in a previous tile.
- \vec{m} has not been defined earlier.

Challenges and contributions

General principle for Load sets

Load a data indexed by \vec{m} just before a tile indexed by \vec{T} if:

- \vec{m} is live-in for \vec{T} , i.e., read but not written **earlier** in \vec{T} .
- \vec{m} has not been loaded in a **previous** tile.
- \vec{m} has not been defined **earlier**.


Tiling defines a schedule on tile+iteration indices, thus “previous” and “earlier”.  **This schedule is not affine in terms of tile sizes.**

Challenges and contributions

General principle for Load sets

Load a data indexed by \vec{m} just before a tile indexed by \vec{T} if:

- \vec{m} is live-in for \vec{T} , i.e., read but not written earlier in \vec{T} .
- \vec{m} has not been loaded in a previous tile.
- \vec{m} has not been defined earlier.

Tiling defines a schedule on tile+iteration indices, thus “previous” and “earlier”.  This schedule is not affine in terms of tile sizes.

Exact case

Reads/writes are functions of iteration points. Can we express the relation “happens before” among iterations in a quasi-affine way?


☛ Yes. Parametric tiling with exact inter-tile reuse is feasible.

Challenges and contributions

General principle for Load sets

Load a data indexed by \vec{m} just before a tile indexed by \vec{T} if:

- \vec{m} is live-in for \vec{T} , i.e., read but not written earlier in \vec{T} .
- \vec{m} has not been loaded in a previous tile.
- \vec{m} has not been defined earlier.

Tiling defines a schedule on tile+iteration indices, thus “previous” and “earlier”.  This schedule is not affine in terms of tile sizes.

Exact case

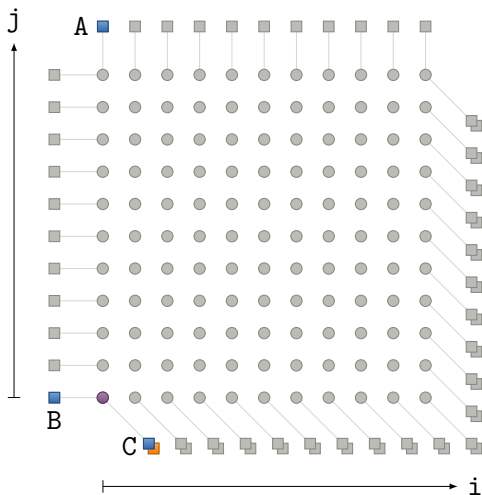
Reads/writes are functions of iteration points. Can we express the relation “happens before” among iterations in a quasi-affine way?

☛ Yes. Parametric tiling with exact inter-tile reuse is feasible.

Approximations

What if contributions of reads/writes are summarized at tile level? Approximated? ☛ No information loss if approximations are “pointwise”. More approximations needed otherwise.

Reads, writes, schedule



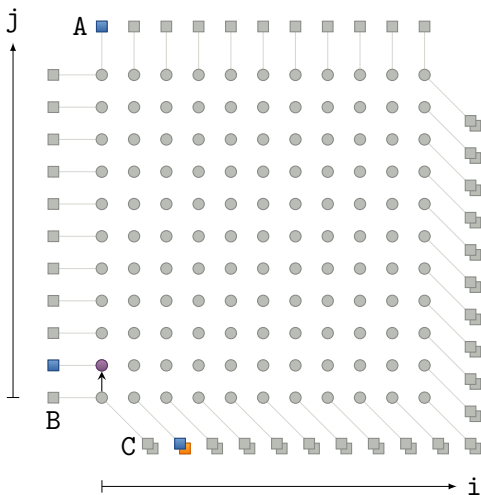
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Reads, writes, schedule



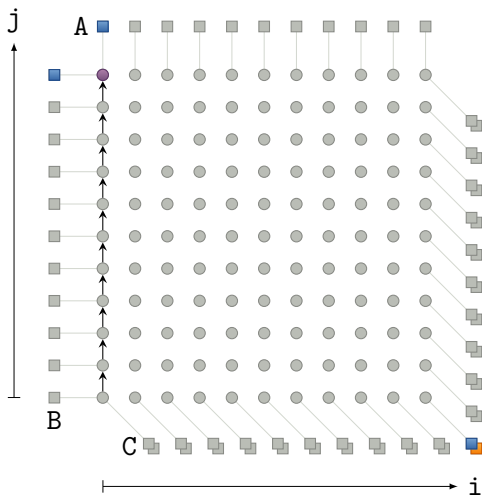
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Reads, writes, schedule



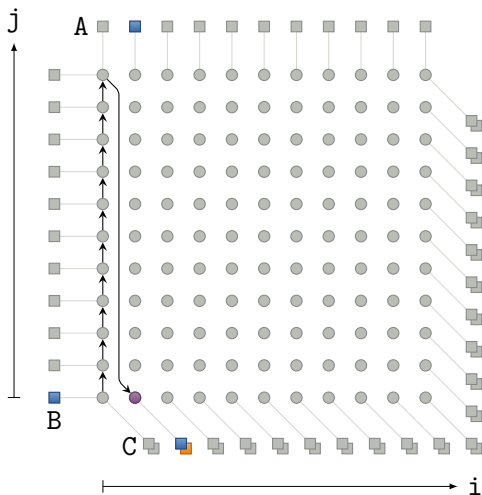
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Reads, writes, schedule



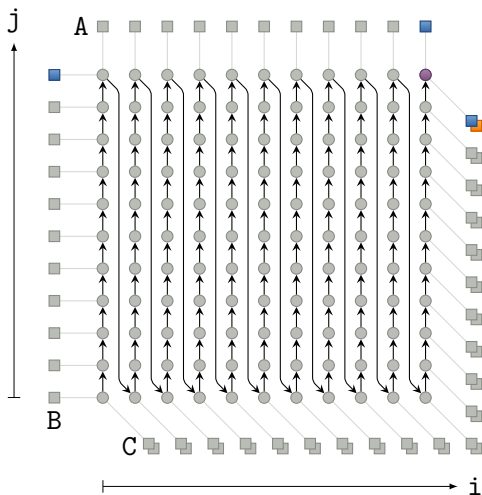
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```


Reads, writes, schedule



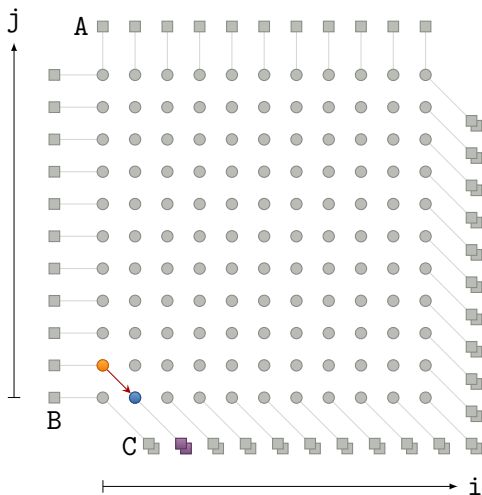
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Dependences



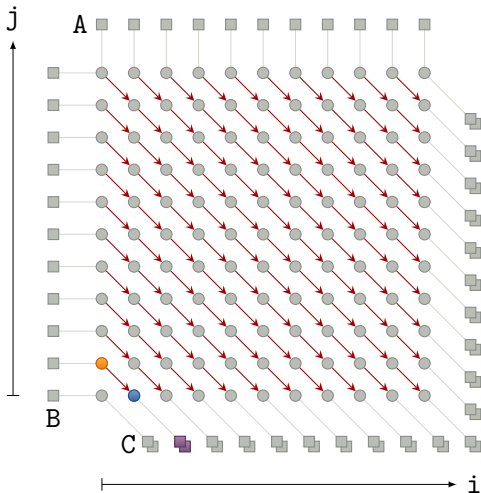
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Dependences



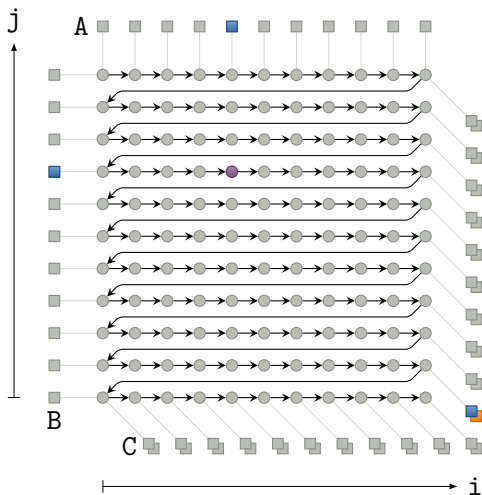
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Scheduling alternatives: loop reversal+interchange



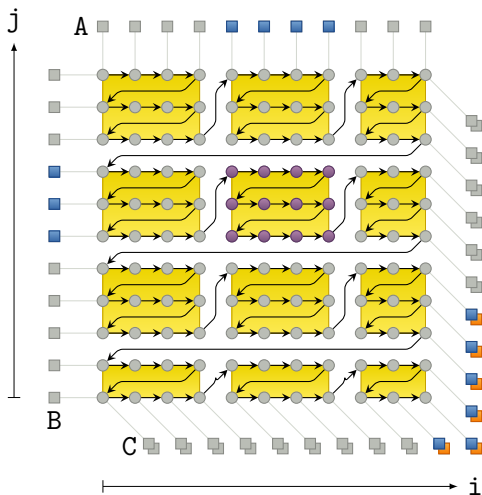
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Scheduling alternatives: loop reversal+interchange+tiling



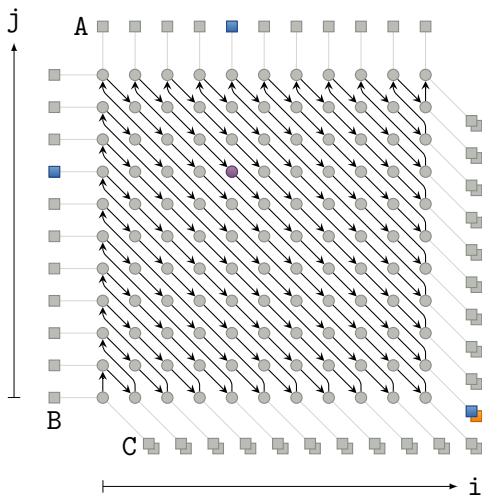
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Scheduling alternatives: loop skewing



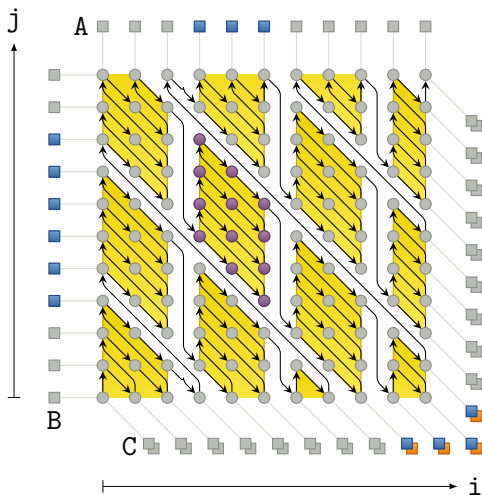
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Scheduling alternatives: loop skewing+tiling



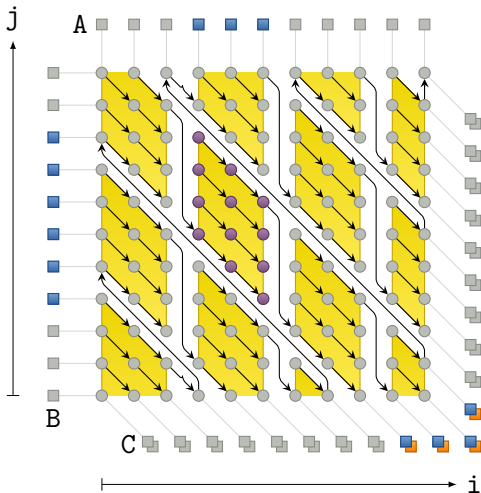
Product of two polynomials:

- arguments in A and B ;
- result in C .

```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

Scheduling alternatives: loop skewing+tiling



+ possibility of **intra-tile parallelism**.

Product of two polynomials:

- arguments in A and B ;
- result in C .

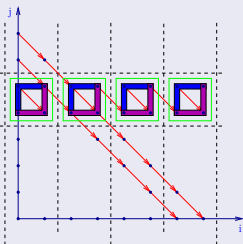
```
for(int k=0; k<2*n-1; k++) {
    C[k] = 0; // S0
}

for(int i=0; i<n; i++) {
    for(int j=0; j<n; j++) {
        C[i+j] += A[i]*B[j]; // S1
    }
}
```

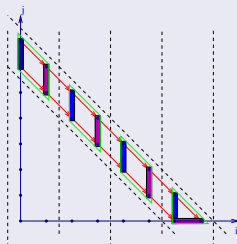

Inter-tile data reuse in a tile strip

```
for(i=0; i<n; i++)
  for(j=0; j<n; j++)
    C[i+j] = C[i+j] + A[i]*B[j];
```

$$(i, j) \mapsto (n - j - 1, i)$$



$$(i, j) \mapsto (i + j, i)$$

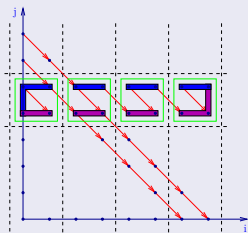


In a tile, **Load** \simeq first read, **Store** \simeq last write.

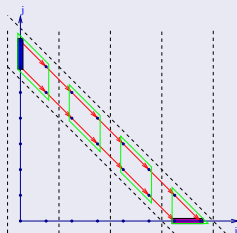
Inter-tile data reuse in a tile strip

```
for(i=0; i<n; i++)
  for(j=0; j<n; j++)
    C[i+j] = C[i+j] + A[i]*B[j];
```

$$(i, j) \mapsto (n - j - 1, i)$$



$$(i, j) \mapsto (i + j, i)$$

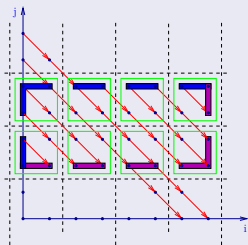


In a tile strip, **Load** \simeq first read, **Store** \simeq last write.

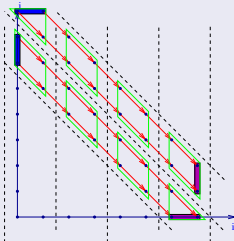
Inter-tile data reuse in a tile strip

```
for(i=0; i<n; i++)  
  for(j=0; j<n; j++)  
    C[i+j] = C[i+j] + A[i]*B[j];
```

$$(i, j) \mapsto (n - j - 1, i)$$

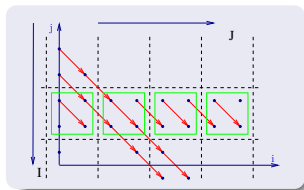


$$(i, j) \mapsto (i + j, i)$$



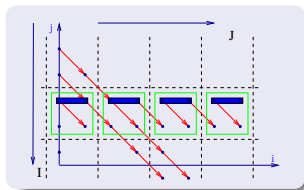
In a **reuse domain**, **Load** \simeq first read, **Store** \simeq last write.
Can actually be adapted to any **parameterized reuse domain**.

Objective: data transfers



- Bound n , tiles of size $b \times b$.
- Tiling with $(i, j) \mapsto (i', j') = (n - j - 1, i)$.
- Access functions $m = i + j = j' + n - i' - 1$.
- Tile origin (I, J) .
- Transfers $\text{Load}_A, \text{Load}_B, \text{Load}_C, \text{Store}_C$.

Objective: data transfers

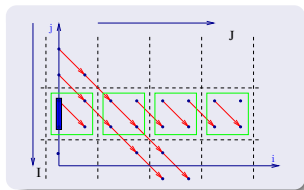


- Bound n , tiles of size $b \times b$.
- Tiling with $(i, j) \mapsto (i', j') = (n - j - 1, i)$.
- Access functions $m = i + j = j' + n - i' - 1$.
- Tile origin (l, J) .
- Transfers $\text{Load}_A, \text{Load}_B, \text{Load}_C, \text{Store}_C$.

Load sets.

$$\text{Load}_A = \{m \mid 0 \leq m \leq n - 1, J \leq m \leq J + b - 1\}$$

Objective: data transfers



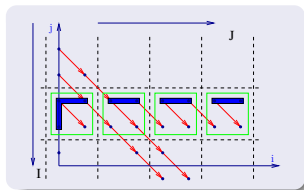
- Bound n , tiles of size $b \times b$.
- Tiling with $(i, j) \mapsto (i', j') = (n - j - 1, i)$.
- Access functions $m = i + j = j' + n - i' - 1$.
- Tile origin (l, J) .
- Transfers $\text{Load}_A, \text{Load}_B, \text{Load}_C, \text{Store}_C$.

Load sets.

$$\text{Load}_A = \{m \mid 0 \leq m \leq n - 1, J \leq m \leq J + b - 1\}$$

$$\text{Load}_B = \{m \mid J = 0, 0 \leq m \leq n - 1, n - l - b \leq m \leq n - l - 1\}$$

Objective: data transfers



- Bound n , tiles of size $b \times b$.
- Tiling with $(i, j) \mapsto (i', j') = (n - j - 1, i)$.
- Access functions $m = i + j = j' + n - i' - 1$.
- Tile origin (l, J) .
- Transfers $\text{Load}_A, \text{Load}_B, \text{Load}_C, \text{Store}_C$.

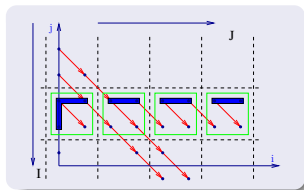
Load sets.

$$\text{Load}_A = \{m \mid 0 \leq m \leq n - 1, J \leq m \leq J + b - 1\}$$

$$\text{Load}_B = \{m \mid J = 0, 0 \leq m \leq n - 1, n - l - b \leq m \leq n - l - 1\}$$

$$\begin{aligned} \text{Load}_C &= \{m \mid 0 \leq m, n - l - b \leq m \leq n - 1 - l, J = 0\} \\ &\cup \{m \mid \max(1, J) \leq m + l - n + 1 \leq \min(n - 1, J + b - 1)\} \end{aligned}$$

Objective: data transfers and local memory sizes



- Bound n , tiles of size $b \times b$.
- Tiling with $(i, j) \mapsto (i', j') = (n - j - 1, i)$.
- Access functions $m = i + j = j' + n - i' - 1$.
- Tile origin (l, J) .
- Transfers $\text{Load}_A, \text{Load}_B, \text{Load}_C, \text{Store}_C$.

Load sets. Local memory sizes with “double-buffering”.

$$\text{Load}_A = \{m \mid 0 \leq m \leq n - 1, J \leq m \leq J + b - 1\}$$

- size $2b$, when $n \geq 2b + 1$: **at least 2 tiles available**.
- size n when $n \leq 2b$: **less than 2 tiles**.

$$\text{Load}_B = \{m \mid J = 0, 0 \leq m \leq n - 1, n - l - b \leq m \leq n - l - 1\}$$

- size b when $n \geq b$: **1 full tile**.
- size n when $n \leq b - 1$: **1 partial tile**.

$$\text{Load}_C = \{m \mid 0 \leq m, n - l - b \leq m \leq n - 1 - l, J = 0\} \\ \cup \{m \mid \max(1, J) \leq m + l - n + 1 \leq \min(n - 1, J + b - 1)\}$$

- size $3b - 1 = (2b - 1) + b$ si $n \geq 2b + 1$: **2 full tiles**.
- size $b + n - 1 = (2b - 1) + (n - b)$ si $b \leq n \leq 2b$: **1 full tile, 1 partial tile**.
- size $2n - 1$ si $n \leq b - 1$: **1 partial tile**.

Outline

- 1 Motivation and challenges
- 2 Parametric analysis
 - Tile index vs tile origin index
 - Exact inter-tile reuse
 - Approximated inter-tile reuse
- 3 Current implementation and results

Tiling, tiles, and schedules

With indices of tiles (tile sizes defined by $\vec{s} = (s_1, \dots, s_n)$)

$$\vec{i} \in \text{Tile}(\vec{T}) \Leftrightarrow \begin{cases} s_1 T_1 \leq i_1 < s_1(T_1 + 1) \\ \vdots \\ s_n T_n \leq i_n < s_n(T_n + 1) \end{cases}$$

- Schedule on iteration points: $\vec{i}' < \vec{i} \Leftrightarrow (\vec{T}', \vec{i}') <_{\text{lex}} (\vec{T}, \vec{i})$.

Tiling, tiles, and schedules

With indices of tiles (tile sizes defined by $\vec{s} = (s_1, \dots, s_n)$)

$$\vec{i} \in \text{Tile}(\vec{T}) \Leftrightarrow \begin{cases} s_1 T_1 \leq i_1 < s_1(T_1 + 1) \\ \vdots \\ s_n T_n \leq i_n < s_n(T_n + 1) \end{cases}$$

- Schedule on iteration points: $\vec{i}' < \vec{i} \Leftrightarrow (\vec{T}', \vec{i}') <_{\text{lex}} (\vec{T}, \vec{i})$.

With indices of tile origins

$$\vec{i} \in \text{Tile}(\vec{l}) \Leftrightarrow \begin{cases} l_1 \leq i_1 < l_1 + s_1 \\ \vdots \\ l_n \leq i_n < l_n + s_n \end{cases} \quad \begin{array}{l} \text{with } \vec{l}, \text{ origin of Tile}(\vec{T}), \\ \text{i.e., } \vec{l} = (s_1 T_1, \dots, s_n T_n). \end{array}$$

- Schedule on iteration points, for a tiling specified by a given tile:

$$\vec{i}' <_{\vec{l}} \vec{i} \Leftrightarrow \vec{i}' <_{\vec{l}} \vec{i} \Leftrightarrow (\vec{l}', \vec{i}') <_{\text{lex}} (\vec{l}, \vec{i}) \text{ and } \vec{l}' \stackrel{\vec{s}}{\equiv} \vec{l}$$

Intuitive expression of Load/Store sets

For $\text{Tile}(\vec{l})$ with data reuse in ReuseDomain:

$$\text{Load}(\vec{l}) = \bigcup_{\vec{i} \in \text{Tile}(\vec{l})} \left(\text{read}(\vec{i}) \setminus \bigcup_{\substack{\vec{i}' < \vec{i} \\ \vec{i}' \in \text{ReuseDomain}}} \text{read}(\vec{i}') \cup \text{write}(\vec{i}') \right)$$

$$\text{Store}(\vec{l}) = \bigcup_{\vec{i} \in \text{Tile}(\vec{l})} \left(\text{write}(\vec{i}) \setminus \bigcup_{\substack{\vec{i}' > \vec{i} \\ \vec{i}' \in \text{ReuseDomain}}} \text{write}(\vec{i}') \right)$$

where $\vec{i}' < \vec{i}$ means that i' is executed before i in the tiled schedule.

Intuitive expression of Load/Store sets

For $\text{Tile}(\vec{i})$ with data reuse in ReuseDomain:

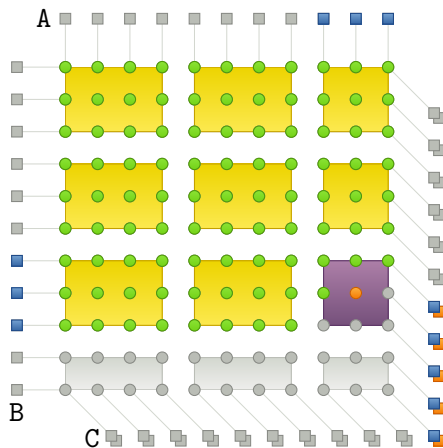
$$\text{Load}(\vec{i}) = \bigcup_{\vec{i}' \in \text{Tile}(\vec{i})} \left(\text{read}(\vec{i}) \setminus \bigcup_{\substack{\vec{i}' < \vec{i} \\ \vec{i}' \in \text{ReuseDomain}}} \text{read}(\vec{i}') \cup \text{write}(\vec{i}') \right)$$

$$\text{Store}(\vec{i}) = \bigcup_{\vec{i}' \in \text{Tile}(\vec{i})} \left(\text{write}(\vec{i}) \setminus \bigcup_{\substack{\vec{i}' > \vec{i} \\ \vec{i}' \in \text{ReuseDomain}}} \text{write}(\vec{i}') \right)$$

where $\vec{i}' < \vec{i}$ means that i' is executed before i in the tiled schedule.

☛ Can we express $\vec{i}' < \vec{i}$ ("happens before") in a parametric way?

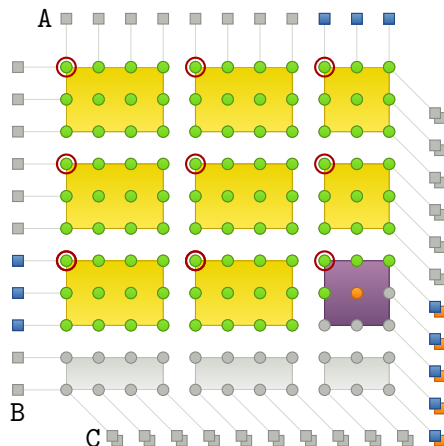
Tiling, relation “happens before” and unaligned tiles



$\vec{i}' < \vec{i}$ iff

- $\vec{i}' \in \text{Tile}(\vec{T})$ and $\vec{i}' \in \text{Tile}(\vec{T}')$
- $(\vec{T}', \vec{i}') <_{\text{lex}} (\vec{T}, \vec{i})$

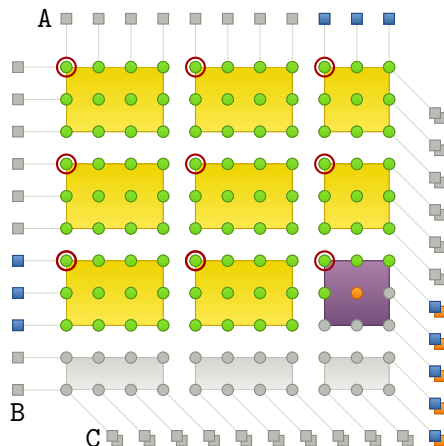
Tiling, relation “happens before” and unaligned tiles



$$\vec{i}' <_{\vec{T}} \vec{i} \text{ iff}$$

- $\vec{i} \in \text{Tile}(\vec{T})$ and $\vec{i}' \in \text{Tile}(\vec{T}')$
- $(\vec{T}', \vec{i}') <_{\text{lex}} (\vec{T}, \vec{i})$ and $\vec{T}' \equiv \vec{T}$

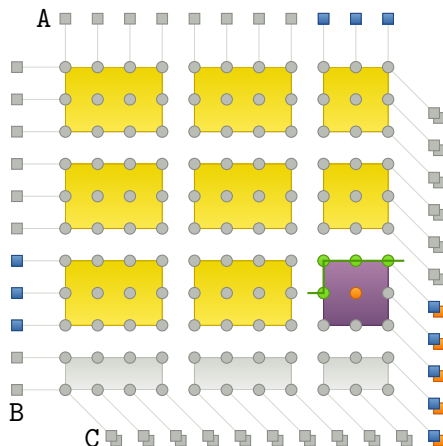
Tiling, relation “happens before” and unaligned tiles



$$\vec{i}' <_{\vec{T}} \vec{i} \text{ iff}$$

- $\vec{i} \in \text{Tile}(\vec{T})$ and $\vec{i}' \in \text{Tile}(\vec{T}')$
 - $\vec{T}' = \vec{T} \wedge \vec{i}' <_{\text{lex}} \vec{i}$
- or
- $$\vec{i}' <_{\text{lex}} \vec{T} \wedge \vec{T}' \stackrel{s}{\equiv} \vec{T} \Leftrightarrow \vec{T} \sqsubset_s \vec{T}'$$

Tiling, relation “happens before” and unaligned tiles



$$\vec{i}' <_{\vec{T}} \vec{i} \text{ iff}$$

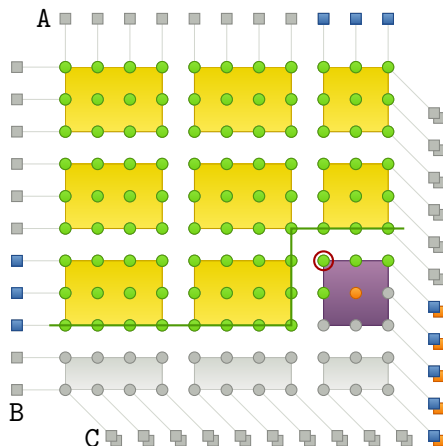
- $\vec{i} \in \text{Tile}(\vec{T})$ and $\vec{i}' \in \text{Tile}(\vec{T}')$

- $\vec{i}' = \vec{i} \wedge \vec{i}' <_{\text{lex}} \vec{i}$

or

$$\vec{i}' <_{\text{lex}} \vec{i} \wedge \vec{i}' \stackrel{s}{\equiv} \vec{i} \Leftrightarrow \vec{i} \sqsubset_s \vec{i}'$$

Tiling, relation “happens before” and unaligned tiles



$$\vec{i}' <_{\vec{T}} \vec{i} \text{ iff}$$

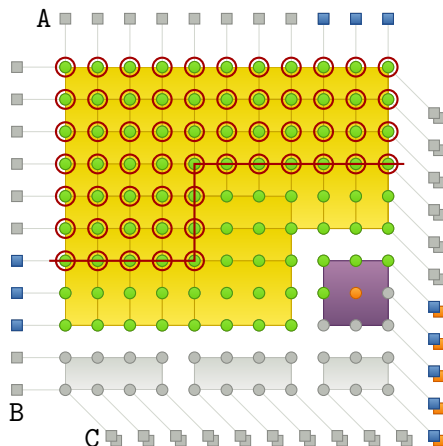
- $\vec{i} \in \text{Tile}(\vec{T})$ and $\vec{i}' \in \text{Tile}(\vec{T}')$

- $\vec{T}' = \vec{T} \wedge \vec{i}' <_{\text{lex}} \vec{i}$

or

$$(i'_1 < h_1) \vee (i'_1 < h_1 + s_1 \wedge i'_2 < h_2)$$

Tiling, relation “happens before” and unaligned tiles



$$\vec{i}' <_{\vec{T}} \vec{i} \text{ iff}$$

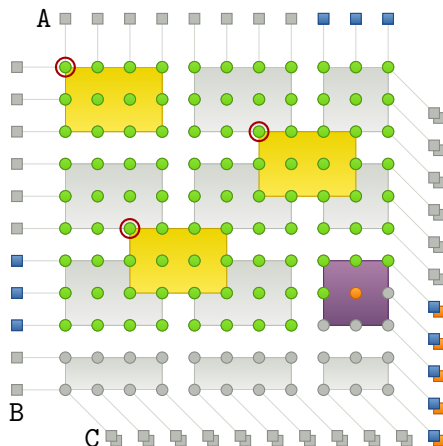
- $\vec{i} \in \text{Tile}(\vec{T})$ and $\vec{i}' \in \text{Tile}(\vec{T}')$

- $\vec{i}' = \vec{T} \wedge \vec{i}' <_{\text{lex}} \vec{i}$

or

$$(l'_1 \leq l_1 - s_1) \vee (l'_1 \leq l_1 \wedge l'_2 \leq l_2 - s_2)$$

Tiling, relation “happens before” and unaligned tiles



$$\vec{i}' <_{\vec{T}} \vec{i} \text{ iff}$$

- $\vec{i} \in \text{Tile}(\vec{T})$ and $\vec{i}' \in \text{Tile}(\vec{T}')$
 - $\vec{T}' = \vec{T} \wedge \vec{i}' <_{\text{lex}} \vec{i}$
- or
- $\vec{i}' <_{\vec{T}'} \vec{i}$: partial order on tiles
 (aligned and unaligned tiles)

Load/Store computations with In/Out sets

Contribution of reads/writes summarized at tile level:

$$\left\{ \begin{array}{l} \text{In}(\vec{T}) = \bigcup_{\vec{i} \in \text{Tile}(\vec{T})} \left(\text{read}(\vec{i}) \setminus \bigcup_{\vec{i}' \in \text{Tile}(\vec{T}), \vec{i}' <_{\text{lex}} \vec{i}} \text{write}(\vec{i}') \right) \\ \text{Out}(\vec{T}) = \bigcup_{\vec{i} \in \text{Tile}(\vec{T})} \text{write}(\vec{i}) \end{array} \right.$$

Load/Store computations with In/Out sets

Contribution of reads/writes summarized at tile level:

$$\left\{ \begin{array}{l} \text{In}(\vec{T}) = \bigcup_{\vec{i} \in \text{Tile}(\vec{T})} \left(\text{read}(\vec{i}) \setminus \bigcup_{\vec{i}' \in \text{Tile}(\vec{T}), \vec{i}' <_{\text{lex}} \vec{i}} \text{write}(\vec{i}') \right) \\ \text{Out}(\vec{T}) = \bigcup_{\vec{i} \in \text{Tile}(\vec{T})} \text{write}(\vec{i}) \end{array} \right.$$

$$\text{Load}(\vec{T}) = \bigcup_{\vec{i} \in \text{Tile}(\vec{T})} \left(\text{read}(\vec{i}) \setminus \bigcup_{\substack{\vec{i}' < \vec{i} \\ \vec{i}' \in \text{ReuseDomain}}} \text{read}(\vec{i}') \cup \text{write}(\vec{i}') \right)$$

$$\bullet \quad \text{Load}(\vec{T}) = \text{In}(\vec{T}) \setminus \left(\bigcup_{\vec{i}' <_{\vec{s}} \vec{i}} \text{In}(\vec{i}') \cup \text{Out}(\vec{i}') \right)$$

Approximations: why?

Some operations *may* execute

- if conditions that are not analyzable.

Some data *may* be accessed

- access functions that are not fully analyzable.

Approximated In/Out sets for tiles ➤ $\overline{\text{In}}$, $\overline{\text{Out}}$, $\underline{\text{Out}}$.

- due to the analysis (e.g., array regions);
- by choice to represent simpler sets (e.g., hyper-rectangles);
- to simplify the analysis (e.g., Fourier-Motzkin).

Approximated Load/Store sets ➤ $\overline{\text{Store}}$, $\overline{\text{Load}}$.

- to simplify code generation;
- to perform communications by blocks;
- to simplify memory allocation;
- ...

Equality of unions

“Exact approximated” load formula

$$\text{Load}(\vec{l}) = \overline{\text{Ra}}_{\vec{l}} \cap ((\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}) \setminus (\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}' \sqsubset_{\vec{s}} \vec{l}))$$

Equality of unions

“Exact approximated” load formula

$$\text{Load}(\vec{l}) = \overline{\text{Ra}}_{\vec{l}} \cap ((\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}) \setminus (\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}' \sqsubset_{\vec{s}} \vec{l}))$$

Simplified “exact” load formula, with **aligned tiles**

$$\text{Load}(\vec{l}) = (\overline{\text{In}} \cup \overline{\text{Out}})(\vec{l}) \setminus (\overline{\text{In}} \cup \overline{\text{Out}})(\vec{l}' \sqsubset_{\vec{s}} \vec{l})$$

Equality of unions

“Exact approximated” load formula

$$\text{Load}(\vec{l}) = \overline{\text{Ra}}_{\vec{l}} \cap ((\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}) \setminus (\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}' \sqsubset_{\vec{s}} \vec{l}))$$

Simplified “exact” load formula, with **aligned tiles**

$$\text{Load}(\vec{l}) = (\overline{\text{In}} \cup \overline{\text{Out}})(\vec{l}) \setminus \bigcup_{\vec{l}' \sqsubset_{\vec{s}} \vec{l}} (\overline{\text{In}} \cup \overline{\text{Out}})(\vec{l}')$$

Equality of unions

“Exact approximated” load formula

$$\text{Load}(\vec{l}) = \overline{\text{Ra}}_{\vec{l}} \cap ((\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}) \setminus (\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}' \sqsubset_{\vec{s}} \vec{l}))$$

Simplified “exact” load formula, with **aligned tiles**

$$\text{Load}(\vec{l}) = F(\vec{l}) \setminus \bigcup_{\vec{l}' \sqsubset_{\vec{s}} \vec{l}} F(\vec{l}')$$

Equality of unions

“Exact approximated” load formula

$$\text{Load}(\vec{l}) = \overline{\text{Ra}}_{\vec{l}} \cap ((\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}) \setminus (\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{l}' \sqsubset_{\vec{s}} \vec{l}))$$

Simplified “exact” load formula, with **aligned tiles** or **all tiles**?

$$\text{Load}(\vec{l}) = F(\vec{l}) \setminus \bigcup_{\vec{l}' \sqsubset_{\vec{s}} \vec{l}} F(\vec{l}') \stackrel{?}{=} F(\vec{l}) \setminus \bigcup_{\vec{l}' \prec_{\vec{s}} \vec{l}} F(\vec{l}')$$

Equality of unions

“Exact approximated” load formula

$$\text{Load}(\vec{T}) = \overline{\text{Ra}}_{\vec{T}} \cap ((\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{T}) \setminus (\overline{\text{In}}' \cup \overline{\text{Out}})(\vec{T}' \sqsubset_{\vec{s}} \vec{T}))$$

Simplified “exact” load formula, with **aligned tiles** or **all tiles**?

$$\text{Load}(\vec{T}) = F(\vec{T}) \setminus \bigcup_{\vec{T}' \sqsubset_{\vec{s}} \vec{T}} F(\vec{T}') \stackrel{?}{=} F(\vec{T}) \setminus \bigcup_{\vec{T}' \prec_{\vec{s}} \vec{T}} F(\vec{T}')$$

Definition (Function stable for unions)

$F : \mathcal{C} \subseteq \mathcal{P}(A) \rightarrow \mathcal{P}(B)$ is *stable for unions* iff $\forall \mathcal{C}', \mathcal{C}'' \subseteq \mathcal{C}$,
 $\bigcup_{X \in \mathcal{C}'} X = \bigcup_{X \in \mathcal{C}''} X \Rightarrow \bigcup_{X \in \mathcal{C}'} F(X) = \bigcup_{X \in \mathcal{C}''} F(X)$.

$$\bigcup_{\vec{T}' \sqsubset_{\vec{s}} \vec{T}} \text{Tile}(\vec{T}') = \bigcup_{\vec{T}' \prec_{\vec{s}} \vec{T}} \text{Tile}(\vec{T}') \stackrel{?}{\Rightarrow} \bigcup_{\vec{T}' \sqsubset_{\vec{s}} \vec{T}} F(\vec{T}') = \bigcup_{\vec{T}' \prec_{\vec{s}} \vec{T}} F(\vec{T}')$$

Pointwise functions

Definition (Function stable for unions)

$F : \mathcal{C} \subseteq \mathcal{P}(\mathcal{A}) \rightarrow \mathcal{P}(\mathcal{B})$ is *stable for unions* iff $\forall \mathcal{C}', \mathcal{C}'' \subseteq \mathcal{C}$,
 $\bigcup_{X \in \mathcal{C}'} X = \bigcup_{X \in \mathcal{C}''} X \Rightarrow \bigcup_{X \in \mathcal{C}'} F(X) = \bigcup_{X \in \mathcal{C}''} F(X)$.

equivalent to

Definition (Pointwise function)

\mathcal{A}, \mathcal{B} two sets, $\mathcal{C} \subseteq \mathcal{P}(\mathcal{A})$. $F : \mathcal{C} \rightarrow \mathcal{P}(\mathcal{B})$ is *pointwise* iff there exists $f : \mathcal{A} \rightarrow \mathcal{P}(\mathcal{B})$ such that $\forall X \in \mathcal{C}$, $F(X) = \bigcup_{x \in X} f(x)$.

Ex: $F(\vec{I}) = (\overline{\text{In}} \cup \overline{\text{Out}})(\vec{I}) = \bigcup_{\vec{i} \in T(\vec{I})} (\overline{\text{read}} \cup \overline{\text{write}})(\vec{i})$.

Pointwise functions

Definition (Function stable for unions)

$F : \mathcal{C} \subseteq \mathcal{P}(A) \rightarrow \mathcal{P}(B)$ is *stable for unions* iff $\forall \mathcal{C}', \mathcal{C}'' \subseteq \mathcal{C}$,
 $\bigcup_{X \in \mathcal{C}'} X = \bigcup_{X \in \mathcal{C}''} X \Rightarrow \bigcup_{X \in \mathcal{C}'} F(X) = \bigcup_{X \in \mathcal{C}''} F(X)$.

equivalent to

Definition (Pointwise function)

A, B two sets, $\mathcal{C} \subseteq \mathcal{P}(A)$. $F : \mathcal{C} \rightarrow \mathcal{P}(B)$ is *pointwise* iff there exists $f : A \rightarrow \mathcal{P}(B)$ such that $\forall X \in \mathcal{C}$, $F(X) = \bigcup_{x \in X} f(x)$.

Ex: $F(\vec{I}) = (\overline{\text{In}} \cup \overline{\text{Out}})(\vec{I}) = \bigcup_{\vec{i} \in T(\vec{I})} (\overline{\text{read}} \cup \overline{\text{write}})(\vec{i})$.

Point-wise approximations

- Largest pointwise under-approximation: $\underline{f}(x) = \bigcap_{Y \in \mathcal{C}, x \in Y} F(Y)$.
- Pointwise over-approximations schemes are possible.

Outline

- 1 Motivation and challenges
- 2 Parametric analysis
- 3 Current implementation and results
 - Current status
 - Script with ISCC
 - Local memory allocation for PolyBench examples

Current implementation and future work

In progress: development of an automated tool

- **ISCC script** (see demo) \Rightarrow complete tool based on ISL.
- Implement approximation schemes: due to code and/or by choice (**complexity issues**). Integrate with PIPS?
- Improve memory size computation: complexity issues, schedules (parallelism), piecewise lattice-based allocation.

To do: experiments with blocking (see also DATE'13)

- **FPGA?** Workstation? GPU? Kalray MPPA?
- Cost model for hierarchical tiling.
- Other schemes of reuse (partial storage).

Pointwise functions

- Useful for other approximations?

Script ISCC 1/3

```
# Inputs
Params := [N, s_1, s_2] -> { : s_1 >= 0 and s_2 >= 0 };
Domain := [N] -> { # Iteration domains
    S_1[k] : 0 <= k < 2N-1;
    S_2[i, j] : 0 <= i, j < N;
} * Params;
Read := [N] -> { # Read access functions
    S_2[i, j] -> A[i];
    S_2[i, j] -> B[j];
    S_2[i, j] -> C[i+j]; } * Domain;
Write := [N] -> { # Write access functions
    S_1[k] -> C[k];
    S_2[i, j] -> C[i+j]; } * Domain;
Theta := [N] -> { # Preliminary mapping
    S_1[k] -> [k, 0, 0];
    S_2[i, j] -> [i+j, i, 1]; };
```

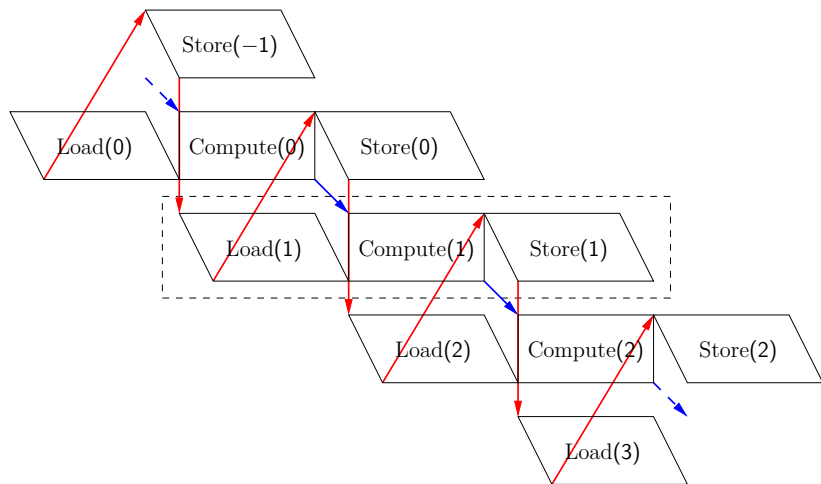
Script ISCC 2/3

```
# Tools for set manipulations
Tiling := [s_1, s_2] -> { # Two dimensional tiling
  [[I_1, I_2] -> [i_1, i_2, k]] -> [i_1, i_2, k] :
    I_1 <= i_1 < I_1 + s_1 and I_2 <= i_2 < I_2 + s_2 };
Coalesce := { [I_1, I_2] -> [[I_1, I_2] -> [i_1, i_2, k]] };
Strip := { [I_1, I_2] -> [I_1, I_2'] };
Prev := { # Lexicographic order
  [[I_1, I_2] -> [i_1, i_2, k]] -> [[I_1, I_2] -> [i_1', i_2', k']] :
    i_1' <= i_1 - 1 or (i_1' <= i_1 and i_2' <= i_2 - 1)
    or (i_1' <= i_1 and i_2' <= i_2 and k' <= k - 1) };
TiledPrev := [s_1, s_2] -> { # Special 'lexicographic' order
  [I_1, I_2] -> [I_1', I_2'] : I_1' <= I_1 - s_1 or
  (I_1' <= I_1 and I_2' <= I_2 - s_2) } * Strip;
TiledNext := TiledPrev^-1;
TiledRead := Tiling.(Theta^-1).Read;
TiledWrite := Tiling.(Theta^-1).Write;
```

Script ISCC 3/3

```
# Set/relation computations
In := Coalesce.(TiledRead - (Prev.TiledWrite));
Out := Coalesce.TiledWrite;
Load := In - ((TiledPrev.In) + (TiledPrev.Out));
Store := Out - (TiledNext.Out);
print coalesce (Load % Params);
print coalesce (Store % Params);
```

Pipelined schedule



Sizes of arrays in local memory

Transformation for tiling	Sequential memory size
jacobi-1d-imper	
$S_0(t, i) \mapsto (t, 2t + i, 0)$	$A[2s_1 + s_2]$
$S_1(t, j) \mapsto (t, 2t + j + 1, 1)$	$B[2s_1 + s_2 - 1]$
jacobi-2d-imper	
$S_0(t, i, j) \mapsto (t, 2t + i, 2t + i + j, 0)$	$A[2s_1 + s_2, \min(2s_1, s_2 + 1) + s_3]$
$S_1(t, i, j) \mapsto (t, 2t + i + 1, 2t + i + j + 1, 1)$	$B[2s_1 + s_2 - 1, \min(2s_1, s_2) + s_3 - 1]$
seidel-2d	
$S_0(t, i, j) \mapsto (t, t + i, 2t + i + j)$	$A \begin{bmatrix} s_1 + s_2 + 1, \\ \min(2s_1 + 2, s_1 + s_2, 2s_2 + 2) + s_3 \end{bmatrix}$
floyd-warshall	
$S_0(k, i, j) \mapsto (k, i, j)$	path $\begin{bmatrix} \max(k + 1, n - k), \\ \max(k + 1, n - k) \end{bmatrix}$

Sizes of arrays in local memory

Transformation for tiling	Pipelined memory size
jacobi-1d-imper	
$S_0(t, i) \mapsto (t, 2t + i, 0)$	$A[2s_1 + 2s_2]$
$S_1(t, j) \mapsto (t, 2t + j + 1, 1)$	$B[2s_1 + 2s_2 - 2]$
jacobi-2d-imper	
$S_0(t, i, j) \mapsto (t, 2t + i, 2t + i + j, 0)$	$A[2s_1 + s_2, \min(2s_1, s_2 + 1) + 2s_3]$
$S_1(t, i, j) \mapsto (t, 2t + i + 1, 2t + i + j + 1, 1)$	$B[2s_1 + s_2 - 1, \min(2s_1, s_2 + 1) + 2s_3 - 2]$
seidel-2d	
$S_0(t, i, j) \mapsto (t, t + i, 2t + i + j)$	$A \begin{bmatrix} s_1 + s_2 + 1, \\ \min(2s_1 + 2, s_1 + s_2, 2s_2 + 2) + 2s_3 \end{bmatrix}$
floyd-warshall	
$S_0(k, i, j) \mapsto (k, i, j)$	path $\begin{bmatrix} \max(k + 1, n - k), \\ \max(k + 1, n - k, 2s_2) \end{bmatrix}$

Merci

Questions ?