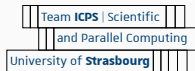# Pipelined Multithreading Generation in a Polyhedral Compiler

January 22nd 2020, IMPACT'20, HiPEAC, Bologna, Italy

Harenome Ranaivoarivony-Razanajato, Cédric Bastoul, Vincent Loechner

University of Strasbourg and Inria Nancy Grand Est

```
1    for (int i = 1; i < N; ++i)
2      A[i] = f1(A[i], A[i - 1]); // S1

3    for (int i = 1; i < N; ++i)
4      B[i] = f2(A[i], B[i - 1]); // S2

5    /* ... */

6    for (int i = 1; i < N; ++i)
7      F[i] = f6(E[i], F[i - 1]); // S6
```

**(a)** Sequential Program

```
1    for (int i = 1; i < N; ++i)
2      A[i] = f1(A[i], A[i − 1]); // S1

3    for (int i = 1; i < N; ++i)
4      B[i] = f2(A[i], B[i − 1]); // S2

5    /* ... */

6    for (int i = 1; i < N; ++i)
7      F[i] = f6(E[i], F[i − 1]); // S6
```



**(a)** Sequential Program          **(b)** Dependency Graph

# Motivating Example

S1(1), thread 1

```
1   for (int i = 1; i < N; ++i)
2     A[i] = f1(A[i], A[i − 1]); // S1

3   for (int i = 1; i < N; ++i)
4     B[i] = f2(A[i], B[i − 1]); // S2

5   /* ... */

6   for (int i = 1; i < N; ++i)
7     F[i] = f6(E[i], F[i − 1]); // S6
```

**(a)** Sequential Program          **(b)** Pipelined Execution

# Motivating Example

S1(1), thread 1 ← S2(1), thread 1

```
1   for (int i = 1; i < N; ++i)
2     A[i] = f1(A[i], A[i - 1]); // S1

3   for (int i = 1; i < N; ++i)
4     B[i] = f2(A[i], B[i - 1]); // S2

5   /* ... */

6   for (int i = 1; i < N; ++i)
7     F[i] = f6(E[i], F[i - 1]); // S6
```

**(a)** Sequential Program                    **(b)** Pipelined Execution
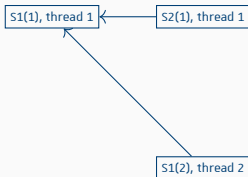
```
1   for (int i = 1; i < N; ++i)
2     A[i] = f1(A[i], A[i - 1]); // S1

3   for (int i = 1; i < N; ++i)
4     B[i] = f2(A[i], B[i - 1]); // S2

5   /* ... */

6   for (int i = 1; i < N; ++i)
7     F[i] = f6(E[i], F[i - 1]); // S6
```



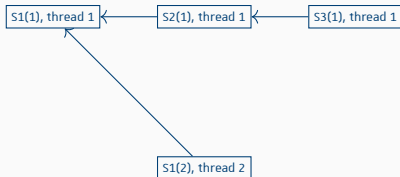**(a)** Sequential Program          **(b)** Pipelined Execution

# Motivating Example

```
1    for (int i = 1; i < N; ++i)
2      A[i] = f1(A[i], A[i − 1]); // S1

3    for (int i = 1; i < N; ++i)
4      B[i] = f2(A[i], B[i − 1]); // S2

5    /* ... */

6    for (int i = 1; i < N; ++i)
7      F[i] = f6(E[i], F[i − 1]); // S6
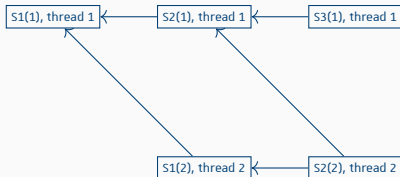```

**(a)** Sequential Program



**(b)** Pipelined Execution

# Motivating Example

```
1    for (int i = 1; i < N; ++i)
2      A[i] = f1(A[i], A[i − 1]); // S1

3    for (int i = 1; i < N; ++i)
4      B[i] = f2(A[i], B[i − 1]); // S2

5    /* ... */

6    for (int i = 1; i < N; ++i)
7      F[i] = f6(E[i], F[i − 1]); // S6
```



**(a)** Sequential Program

**(b)** Pipelined Execution

```
1    for (int i = 1; i < N; ++i)
2      A[i] = f1(A[i], A[i - 1]); // S1

3    for (int i = 1; i < N; ++i)
4      B[i] = f2(A[i], B[i - 1]); // S2

5    /* ... */

6    for (int i = 1; i < N; ++i)
7      F[i] = f6(E[i], F[i - 1]); // S6
```
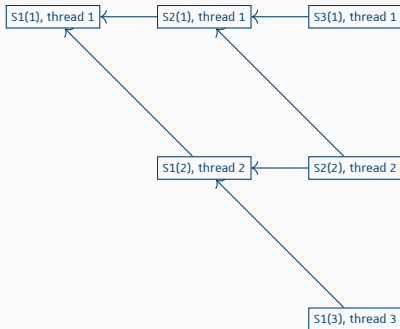
**(a)** Sequential Program

**(b)** Pipelined Execution

# Motivating Example

```
1   for (int i = 1; i < N; ++i)
2     A[i] = f1(A[i], A[i - 1]); // S1

3   for (int i = 1; i < N; ++i)
4     B[i] = f2(A[i], B[i - 1]); // S2

5   /* ... */

6   for (int i = 1; i < N; ++i)
7     F[i] = f6(E[i], F[i - 1]); // S6
```

**(a)** Sequential Program

```
1    #pragma omp parallel
2    {
3        #pragma omp for schedule(static) ordered nowait
4        for (int i = 1; i < N; ++i)
5          #pragma omp ordered
6          A[i] = f1(A[i], A[i - 1]); // S1
7        #pragma omp for schedule(static) ordered nowait
8        for (int i = 1; i < N; ++i)
9          #pragma omp ordered
10         B[i] = f2(A[i], B[i - 1]); // S2

11       /* ... */

12       #pragma omp for schedule(static) ordered nowait
13       for (int i = 1; i < N; ++i)
14         #pragma omp ordered
15         F[i] = f6(E[i], F[i - 1]); // S6
16   }
```

**(b)** Pipelined OpenMP target program

# Motivating Example

```
1   for (int i = 1; i < N; ++i)
2     A[i] = f1(A[i], A[i - 1]); // S1

3   for (int i = 1; i < N; ++i)
4     B[i] = f2(A[i], B[i - 1]); // S2

5   /* ... */

6   for (int i = 1; i < N; ++i)
7     F[i] = f6(E[i], F[i - 1]); // S6
```

**(a)** Sequential Program

```
1    #pragma omp parallel
2    {
3      #pragma omp for schedule(static) ordered nowait
4      for (int i = 1; i < N; ++i)
5        #pragma omp ordered
6        A[i] = f1(A[i], A[i - 1]); // S1
7      #pragma omp for schedule(static) ordered nowait
8      for (int i = 1; i < N; ++i)
9        #pragma omp ordered
10       B[i] = f2(A[i], B[i - 1]); // S2

11     /* ... */

12     #pragma omp for schedule(static) ordered nowait
13     for (int i = 1; i < N; ++i)
14       #pragma omp ordered
15       F[i] = f6(E[i], F[i - 1]); // S6
16   }
```

**(b)** Pipelined OpenMP target program

Speedup: 2.89

6 stages on an Intel Xeon E5-2620v3 @ 2.40 GHz, with $N = 100,000$

- Identifying software pipelines in a polyhedral compiler

- Generate pipelined multithreading using OpenMP

**Pipelined Multithreading Generation in a Polyhedral Compiler**, Harenome Razanajato et al.

2

# Polyhedral Model

- `#pragma` based API for shared memory parallelism
- Worksharing constructs
  - `#pragma omp for`
  - `#pragma omp task`

**Pipelined Multithreading Generation in a Polyhedral Compiler**, Harenome Razanajato et al.

3

# OpenMP

- `#pragma` based API for shared memory parallelism
- Worksharing constructs
  - `#pragma omp for`
  - `#pragma omp task`
- Synchronization
  - `#pragma omp barrier`: explicit synchronization barrier
  - `omp_set_lock()` and `omp_unset_lock()`: explicit lock mechanism

- `#pragma` based API for shared memory parallelism
- Worksharing constructs
  - `#pragma omp for`
  - `#pragma omp task`
- Synchronization
  - `#pragma omp barrier`: explicit synchronization barrier
  - `omp_set_lock()` and `omp_unset_lock()`: explicit lock mechanism
- Clauses
  - `nowait` clause on worksharing constructs: omit the implicit barrier at the end of a worksharing construct
  - `ordered` clause on worksharing constructs: sequentialize a region

# Polyhedral Model

- Goal: maximize the number of pipeline stages

- Dependence analysis: identify Strongly Connected Components

**Pipelined Multithreading Generation in a Polyhedral Compiler**, Harenome Razanajato et al.

4

# Sequential Loop Fission

```
1   for (int i = 2; i < N; ++i) {
2     a[i] = h[i − 1] + R[i]; // S1
3     b[i] = a[i − 1] + a[i]; // S2
4     c[i] = b[i − 1] + b[i]; // S3
5     d[i] = c[i − 1] + c[i]; // S4
6     e[i] = d[i − 2] + d[i − 1]; // S5
7     f[i] = e[i − 2] + e[i − 1]; // S6
8     g[i] = f[i] + X[i]; // S7
9     h[i] = g[i] + Y[i]; // S8
10    u[i] = v[i − 1] + d[i]; // S9
11    v[i] = u[i] + Z[i]; // S10
12  }
```

**(a)** Original loop body

```
1   for (int i = 2; i < N; ++i) {
2     a[i] = h[i − 1] + R[i]; // S1
3     b[i] = a[i − 1] + a[i]; // S2
4     c[i] = b[i − 1] + b[i]; // S3
5     d[i] = c[i − 1] + c[i]; // S4
6     e[i] = d[i − 2] + d[i − 1]; // S5
7     f[i] = e[i − 2] + e[i − 1]; // S6
8     g[i] = f[i] + X[i]; // S7
9     h[i] = g[i] + Y[i]; // S8
10  }
11  for (int i = 2; i < N; ++i) {
12    u[i] = v[i − 1] + d[i]; // S9
13    v[i] = u[i] + Z[i]; // S10
14  }
```

**(b)** Fission of Strongly Connected Components

The safe use of the `nowait` clause between two **parallel** loops requires that there are no dependencies between the loops or that:

- the sizes of the iteration domains are equal
- the chunk size is either the same or not specified
- both loops are bound to the same parallel region
- none of the loops is associated with a SIMD construct
- the second loop depends only on the same logical iteration of the first loop

# Relaxed conditions on the nowait clause for ordered loops

The safe use of the `nowait` clause between two **ordered** loops requires that there are no dependencies between the loops or that:

- the sizes of the iteration domains are equal
- the chunk size is either the same or not specified
- both loops are bound to the same parallel region
- none of the loops is associated with a SIMD construct
- ~~the second loop depends only on the same logical iteration of the first loop~~

**Pipelined Multithreading Generation in a Polyhedral Compiler**, Harenome Razanajato et al.

7

# Relaxed conditions on the nowait clause for ordered loops

The safe use of the `nowait` clause between two **ordered** loops requires that there are no dependencies between the loops or that:

- the sizes of the iteration domains are equal
- the chunk size is either the same or not specified
- both loops are bound to the same parallel region
- none of the loops is associated with a SIMD construct
- ~~the second loop depends only on the same logical iteration of the first loop~~
- the second loop depends on the same logical iteration **or previous logical iterations** of the first loop

```
1   #pragma omp parallel
2   {
3     #pragma omp for nowait
4     for (int i = 0; i < N; ++i)
5       A[i] = f1(A[i]);
6     #pragma omp for
7     for (int i = 0; i < N; ++i)
8       B[i] = f2(B[i], A[i]);
9   }
```

**(a)** Parallel `for` and `nowait`

```
1    #pragma omp parallel
2    {
3      #pragma omp for ordered nowait
4      for (int i = 0; i < N; ++i)
5        #pragma omp ordered
6        A[i] = f1(A[i]);
7      #pragma omp for ordered
8      for (int i = 0; i < N; ++i)
9        #pragma omp ordered
10       B[i] = f2(B[i], A[i-1]);
11   }
```

**(b)** Ordered `for` and `nowait`

- Annotate sequential loops with `#pragma omp for ordered`

- Enclose sequential loop bodies in `#pragma omp ordered` regions

- Annotate loops with `nowait` where possible

- Optimize by reverting `ordered` loops without `nowait` clauses to `#pragma omp single` regions

# Polyhedral Model

- Loop blocking and loop fusion

- `#pragma omp for schedule(static, 1)` on the blocking loop

- `omp_set_lock()` and `omp_unset_lock()` before and after each loop of the pipeline

- up to $n \times m$ locks required for $m$ pipeline stages over $n$ threads

# Explicit synchronization

```
1   #pragma omp parallel
2   {
3       #pragma omp for ordered nowait
4       for (size_t i = 1; i < N; ++i)
5           #pragma omp ordered
6           A[i] = f(A[i], A[i - 1]);

7       /* Other stages */

8       #pragma omp for ordered
9       for (size_t i = 1; i < N; ++i)
10          #pragma omp ordered
11          F[i] = f(E[i], F[i - 1]);
12  }
```

```
1   omp_lock_t** locks;
2   #pragma omp parallel
3   {
4       /* Choose num_threads, block_size, block_count. */
5       /* Allocate, initialize and set the locks. */
6       #pragma omp for schedule(static, 1)
7       for (size_t block = 0; block < block_count; ++block)
8           /* Local loop bounds and indexes. */
9           const size_t start = 1 + block * block_size;
10          const size_t end = MIN(start + block_size, N);
11          const size_t self = block % num_threads;
12          const size_t next = (block + 1) % num_threads;

13          omp_set_lock(&locks[self][0]);
14          for (size_t i = start; i < end; ++i)
15              A[i] = f(A[i], A[i-1]);
16          omp_unset_lock(&locks[next][0]);
17          /* Other stages of the pipeline */
18          omp_set_lock(&locks[self][5]);
19          for (size_t i = start; i < end; ++i)
20              F[i] = f(E[i], F[i-1]);
21          omp_unset_lock(&locks[next][5]);
22      }
23      /* Destroy and free locks. */
24  }
```

**(a)** Original program      **(b)** Pipelined OpenMP target program

**Pipelined Multithreading Generation in a Polyhedral Compiler**, Harenome Razanajato et al.

- Tested on an Intel Xeon E5-2620v3 @ 2.40 GHz, linux 5.3.7

- Code compiled using `gcc 9.2.1` and `clang 9.0.0` with options `-O3 -march=native -fopenmp`

- FIFO scheduling enabled and process priority set to 75

| benchmark | parallel loops | stages |
|---|---|---|
| teaser | 0 | 5 |
| van_dongen[1] | 0 | 2 |
| wdf[2] | 0 | 2 |
| mix | 1 | 2 |

---

[1] (Vincent H Van Dongen, Guang R Gao, and Qi Ning. "A polynomial time method for optimal software pipelining". In: *Parallel Processing: CONPAR 92—VAPP V*. Springer, 1992, pp. 613–624)

[2] (Alfred Fettweis. "Wave digital filters: Theory and practice". In: *Proceedings of the IEEE* 74.2 [1986], pp. 270–327)
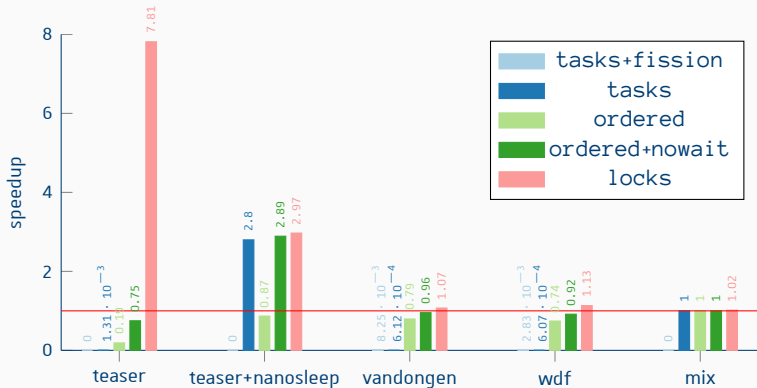
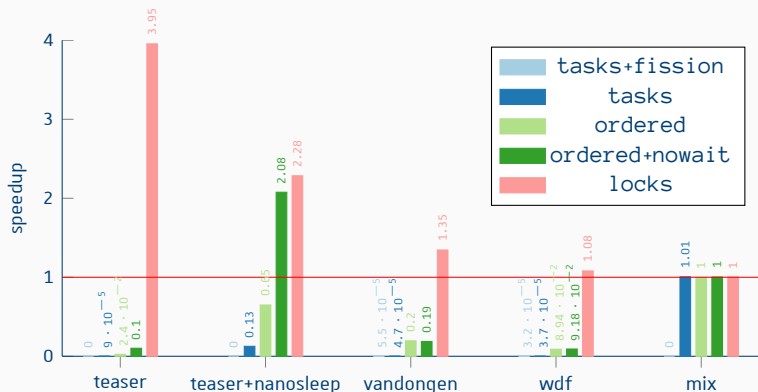**Figure 7:** Speedups or slowdowns over sequential version

**Figure 8:** Speedups or slowdowns over sequential version

- Contributions:
  - Identifying software pipelines in a polyhedral compiler
  - Generating pipelined multithreading

# Contributions and Future Work

- Contributions:
    - Identifying software pipelines in a polyhedral compiler
    - Generating pipelined multithreading
- Future work:
    - Integration in an automatic parallelizer
    - Investigate methods to determine optimal block sizes

Appendix

[1] Alfred Fettweis. "Wave digital filters: Theory and practice". In: *Proceedings of the IEEE* 74.2 (1986), pp. 270–327.

[2] Harenome Razanajato, Cédric Bastoul, and Vincent Loechner. "Pipelined Multithreading Generation in a Polyhedral Compiler". In: *IMPACT 2020, 10th International Workshop on Polyhedral Compilation Techniques*. Bologna, Italy, Jan. 2020.

[3] Vincent H Van Dongen, Guang R Gao, and Qi Ning. "A polynomial time method for optimal software pipelining". In: *Parallel Processing: CONPAR 92—VAPP V*. Springer, 1992, pp. 613–624.