# Comparison of the performance of scientific calculation codes

Lin Sinan

Supervisors: Prof. Sid Tuoati

Université Cote d'Azur
EIT Digital

June 25, 2020

UNIVERSITÉ
CÔTE D'AZUR

# Table of Contents

UNIVERSITÉ
**CÔTE D'AZUR**

# Table of Contents

UNIVERSITÉ
**CÔTE D'AZUR**

Figure: Where's the Real Bottleneck in Scientific Computing?
source: American Scientist

A great amount of execution time of a scientific calculation is spent on loops, such as, matrix multiplication. Modern compilers have been developed to make them faster by using loop optimization techniques:

### Loop Optimization
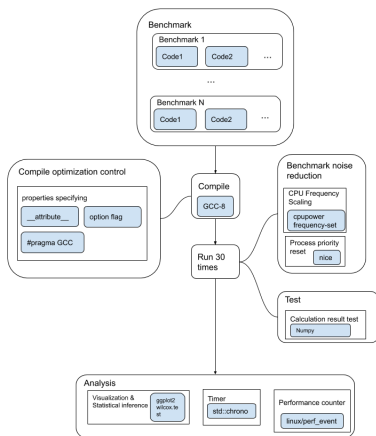
- Loop Unrolling
- Loop Interchange
- Loop Blocking

UNIVERSITÉ
CÔTE D'AZUR

# Procedure



Figure: Benchmarking Procedure

# Table of Contents

UNIVERSITÉ
**CÔTE D'AZUR**

# Loop Unrolling

---

**Algorithm 1** Array Addition

1: $A \leftarrow$ an 32-bit float array with length 100
2: $b \leftarrow$ a 32-bit float scalar
3: i = 0
4: **for** $i < 100$ **do**
5:     $A[i] \leftarrow A[i] + b$
6:     $i \leftarrow i + 1$
7: **EndFor**

---

**Algorithm 1** Array Addition with loop unrolling

1: $A \leftarrow$ an 32-bit float array with length 100
2: $b \leftarrow$ a 32-bit float scalar
3: i = 0
4: **for** $i < 100$ **do**
5:     $A[i] \leftarrow A[i] + b$
6:     $A[i + 1] \leftarrow A[i + 1] + b$
7:     $i \leftarrow i + 2$
8: **EndFor**

---

**Loop unrolling** replicates the body of the loop to reduce loop overhead.

UNIVERSITÉ
**CÔTE D'AZUR**

# Loop Unrolling

**Algorithm 3** Assembly Instructions of Array Addition

1: **procedure**
2:     $r1 \leftarrow$ address of A array
3:     $r2 \leftarrow$ address of A plus offset 400 (we assume A is a 32-bit float array)
4:     $r3 \leftarrow$ address of b
5:     $r4 \leftarrow$ temporary registers
6:     loop:
7:         load r4, 0(r1)
8:         add r4, r4, r3
9:         save r4, 0(r1)
10:        add r1, r1, 4
11:        bne r1, r2, loop

**Algorithm 4** Assembly Instructions with unrolling factor 1

1: **procedure**
2:     $r1 \leftarrow$ address of A array
3:     $r2 \leftarrow$ address of A plus offset 400 (we assume A is a 32bit int array)
4:     $r3 \leftarrow$ address of b
5:     $r4 \leftarrow$ temporary registers
6:     loop:
7:         load r4, 0(r1)
8:         add r4, r4, r3
9:         save r4, 0(r1)
10:        load r4, 4(r1)
11:        add r4, r4, r3
12:        save r4, 4(r1)
13:        add r1, r1, 8
14:        bne r1, r2, loop

Algorithm1 has 100 loops and 5 instructions in the loop body, so the number of instructions is 500. As for Algorithm2, the the number of instructions is 400, since it has 8 instructions in loop body, and only 50 loops.

UNIVERSITÉ
**CÔTE D'AZUR**

# Unrolling

```
// benchmarking starts
    int i,j,k;
    for (i = 0; i < N; i++){
        for (j = 0; j < M; j++){
#ifdef UNROLL_OPTION
#pragma GCC unroll 7
#endif
            for (k = 0; k< P; k++){
                C[i][j] = C[i][j] + A[i][k] * B[k][j];
            }
        }
    }
```

Figure: Common implementation

```
int i,j,k;
for (i = 0; i < N; i++){
    for (j = 0; j < M; j++){
        for (k = 0; k + 7 - 1 < P; k+=7) {
            C[i][j] = C[i][j] + A[i][k] * B[k][j];
            C[i][j] = C[i][j] + A[i][k+1] * B[k+1][j];
            C[i][j] = C[i][j] + A[i][k+2] * B[k+2][j];
            C[i][j] = C[i][j] + A[i][k+3] * B[k+3][j];
            C[i][j] = C[i][j] + A[i][k+4] * B[k+4][j];
            C[i][j] = C[i][j] + A[i][k+5] * B[k+5][j];
            C[i][j] = C[i][j] + A[i][k+6] * B[k+6][j];
        }
        // make sure the remaining loop are used.
        for (; k < P; k++){
            C[i][j] = C[i][j] + A[i][k] * B[k][j];
        }
    }
}
```

Figure: Manually unroll 6 times

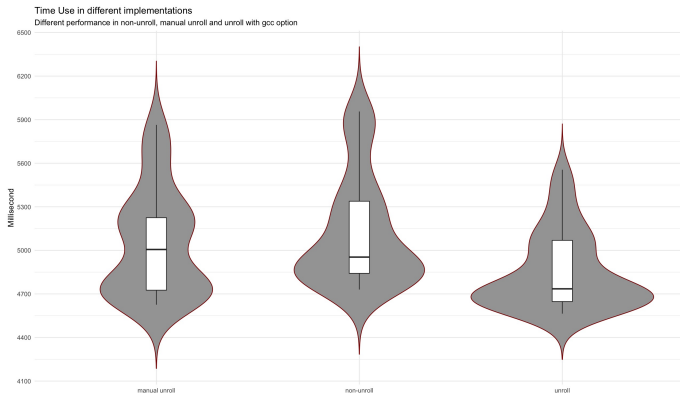Benchmarking code of orginal, manual unrolling and unroll with
-funroll-loops

# Result



Figure: Violin plot of manually unroll, non-unroll and unroll with compiler option

## Result

Table: Unrolling Benchmark Result and Statistical inference

|  | avg time (*ms*) | std | # of instructions (*million*) |
|---|---|---|---|
| unroll | 4848.607 | 262.2337 | 12552.76 |
| no unroll | 5126.94 | 374.7673 | 23645.26 |
| manually unroll | 5049.016 | 365.0488 | 12071.26 |

UNIVERSITÉ
**CÔTE D'AZUR**

# Result

## Paired samples Wilcoxon test result

| Null Hypothesis | P value | Result |
|---|---|---|
| Unroll with compiler option is not faster than non-unroll | <0.0001 | Refuse |
| Unroll with compiler option is not faster than manual-unroll | <0.0001 | Refuse |
| Manual-unroll with compiler option is not faster than non-unroll | 0.1998 | Accept |

UNIVERSITÉ
CÔTE D'AZUR

# Unrolling



Figure: unroll with compiler option

Figure: Manually unroll

Registers are not allocated well for manual version
and increase stalls in instruction pipeline.

UNIVERSITÉ
**CÔTE D'AZUR**

# Cache



Figure: Data transfer among cpu, cache and RAM

When the computer read from or write to a location in main memory, it first checks whether a copy of that data is in the cache. If so, the processor immediately reads from or writes to the cache, which is much faster than reading from or writing to main memory.

# Loop Interchange

**Loop Interchange** : change the order of loop to improve the cache locality.



Figure: Layout of 2D array in memory(row-major)

# Loop Interchange

**Algorithm 7** Matrix multiplication in ijk order
___
1:  int C[N][M];
2:  int A[N][P];
3:  int B[P][M];
4:  for (int i = 0; i < n; i++)
5:      for (int j = 0; j < M; j++)
6:          for (int k = 0; k < P; k++)
7:              C[i][j] += A[i][k] * B[k][j];
___



Figure: Matrix multiplication in ijk order

# Loop Interchange

---

**Algorithm 8** Matrix multiplication in ikj order

---
1: int C[N][M];
2: int A[N][P];
3: int B[P][M];
4: for (int i = 0; i < n; i++)
5:    for (int k = 0; k < P; k++)
6:       for (int j = 0; j < M; j++)
7:          C[i][j] += A[i][k] * B[k][j];
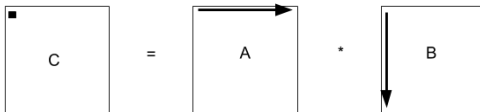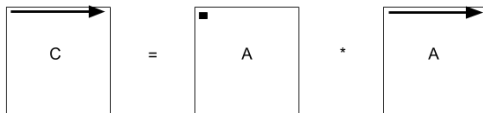
---



Figure: Matrix multiplication in ikjk order

UNIVERSITÉ
CÔTE D'AZUR

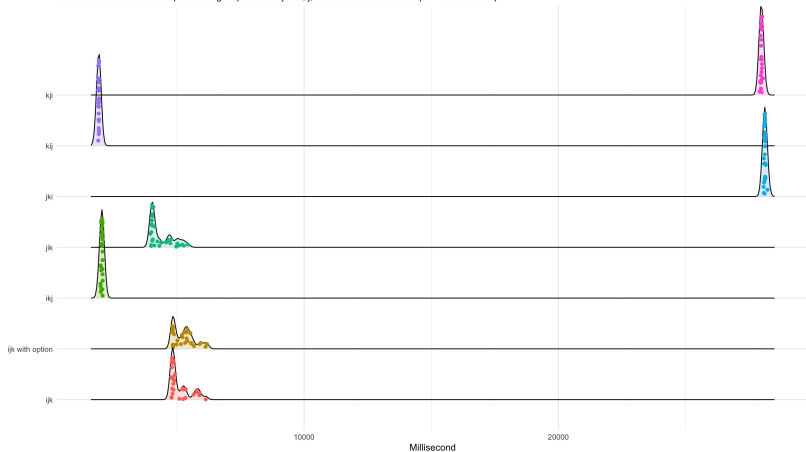Figure: Different loop orders lead to different cache access pattern

Figure: Ridge plot of different loop orders result

Table: Interchange Benchmark Result

|                  | avg time (ms) | std      | LLC cache miss[1] (million) |   |
|------------------|---------------|----------|------------------------------|---|
| ijk with option  | 5250.45       | 411.0194 | 244.04                       |   |
| ijk              | 5118.148      | 374.77   | 229.61                       |   |
| jik              | 4390.417      | 469.84   | 170.75                       |   |
| ikj              | 2044.175      | 29.38    | 136.03                       |   |
| kij              | 1922.871      | 12.56    | 141.40                       |   |
| kji              | 27982.27      | 25.30    | 364.27                       |   |
| jki              | 28122.32      | 30.75    | 301.49                       |   |

UNIVERSITÉ
CÔTE D'AZUR

```cpp
int i,j,k;
for (i = 0; i < N; i++){
    for (j = 0; j < M; j++){
        for (k = 0; k< P; k++){
            C[i][j] = C[i][j] + A[i][k] * B[k][j];
        }
    }
}
```

Figure: Common implementation

```cpp
int i0,j0,k0,i,j,k;
for (i0 = 0; i0 < N; i0 += BLOCK_SIZE)
    for (k0 = 0; k0 < P; k0 += BLOCK_SIZE)
        for (j0 = 0; j0 < M; j0 += BLOCK_SIZE)
            for (i = i0; i < std::min(i0 + BLOCK_SIZE, N); i++)
                for (k = k0; k < std::min(k0 + BLOCK_SIZE, P); k++)
                    for (j = j0; j < std::min(j0 + BLOCK_SIZE, M); j++)
                        C[i][j] = C[i][j] + A[i][k] * B[k][j];
```

Figure: Matrix multiplication with Blocking

UNIVERSITÉ
CÔTE D'AZUR

# Loop Blocking
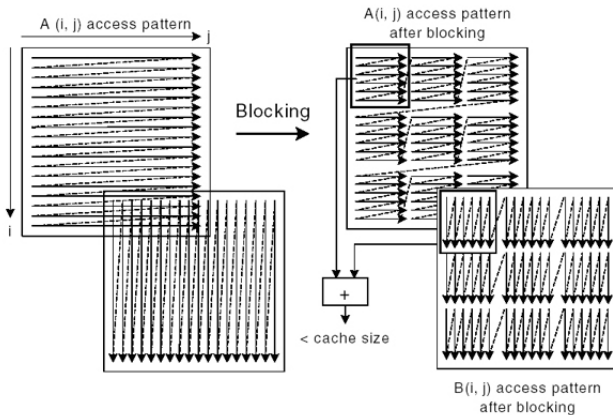

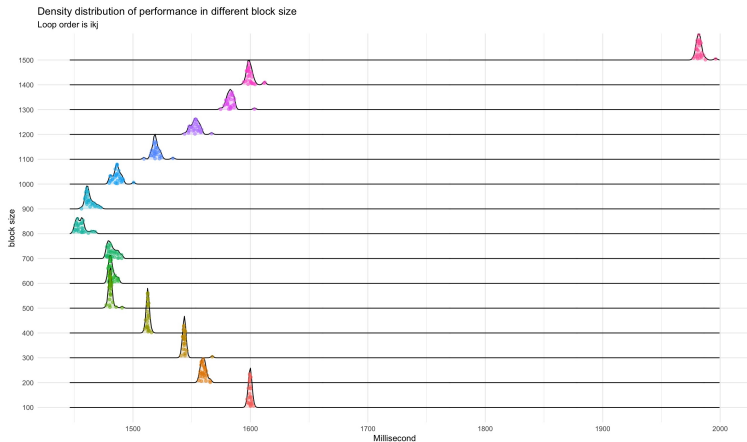
Figure: Loop Blocking visualization(from intel corp)

# result



Figure: Ridge plot of performance in different block size

| Block Size | avg time (ms) | LLC Cache miss (million) | Branch miss (million) |
|---|---|---|---|
| 100 | 1599.766 | 4.01 | 34.43 |
| 300 | 1544.466 | 9.97 | 11.29 |
| 500 | 1481.069 | 4.05 | 6.76 |
| 700 | 1482.241 | 5.85 | 6.76 |
| 900 | 1463.007 | 12.66 | 4.51 |
| 1100 | 1519.416 | 44.01 | 4.51 |
| 1300 | 1583.119 | 90.78 | 4.51 |
| 1500 | 1982.382 | 135.07 | 2.25 |

Figure: Loop Blocking result)

# Table of Contents

UNIVERSITÉ
**CÔTE D'AZUR**

# summary

**SUMMARY**

- Compilers can help to optimize your code well.
- Compilers cannot give an optimal solution sometimes.
- The property of cache is important for program performance.

UNIVERSITÉ
**CÔTE D'AZUR**