# Optimization techniques

## Optimization training at CINES

Adrien CASSAGNE

adrien.cassagne@inria.fr

2015/12/01

# Contents

# Contents

# Pre-compute things when it is possible

- We can divide a computational code in 3 parts
  1. Pre-processing: allocations, initializations, reading of inputs
  2. Solver or massively computational part
  3. Post-processing: writing of outputs, deallocations
- In many cases the solver part is the most time consuming
- When it is possible, we have to pre-compute things in the pre-processing part
  - The idea is to reduce the total number of operations in the code

# Pre-processing example

```
1  void main()
2  {
3    // pre-processing part
4    int n = 1000;
5    float *A, *B, *C, *D;
6    A = new float[n]; // in
7    B = new float[n]; // in
8    C = new float[n]; // in
9    D = new float[n]; // out
10   randomInit(A, n); randomInit(B, n); randomInit(C, n);
11   for(int i = 0; i < n; i++)
12     D[i] = 0;
13
14   // solver or computational part
15   for(int j = 0; j < n; j++)
16     for(int i = 0; i < n; i++)
17       D[i] = D[i] + (A[i] + B[i]) * C[j];
18
19   // post-processing part
20   delete[] A;
21   delete[] B;
22   delete[] C;
23   delete[] D;
24 }
```

Simple computational code, not optimized

- The total number of operations (flops) is $n \times n \times 3$

# Pre-processing example

```
1   void main()
2   {
3     // pre-processing part
4     int n = 1000;
5     float *A, *B, *C, *D, *AB;
6     A  = new float[n]; // in
7     B  = new float[n]; // in
8     AB = new float[n]; // pre-processing
9     C  = new float[n]; // in
10    D  = new float[n]; // out
11    randomInit(A, n); randomInit(B, n); randomInit(C, n);
12    for(int i = 0; i < n; i++) {
13      D [i] = 0;
14      AB[i] = A[i] + B[i]; // we pre-compute A + B
15    }
16
17    // solver or computational part
18    for(int j = 0; j < n; j++)
19      for(int i = 0; i < n; i++)
20        D[i] = D[i] + AB[i] * C[j];
21
22    // post-processing part ...
23  }
```

Simple computational, code optimized

- The total number of operations (flops) is $n + (n \times n \times 2)$
- Can we do better than that?

# Pre-processing example

```
1   void main()
2   {
3     // pre-processing part
4     int n = 1000;
5     float *A, *B, *C, *D, *AB;
6     A  = new float[n]; // in
7     B  = new float[n]; // in
8     C  = new float[n]; // in
9     D  = new float[n]; // out
10    randomInit(A, n); randomInit(B, n); randomInit(C, n);
11    float sumC = 0;
12    for(int i = 0; i < n; i++)
13      sumC += C[i]; // we pre-compute the sum of C[0] + C[1] + ... + C[n -1]
14
15    // solver or computational part
16    for(int i = 0; i < n; i++)
17      D[i] = (A[i] + B[i]) * sumC;
18
19    // post-processing part ...
20  }
```

Simple computational code, well optimized

- The total number of operations (flops) is $3n$

## Branch instructions

- Branch instructions (alias `if`, `switch`, etc) create bubbles in the processor pipeline
- The pipeline can't be fully filled
- We have to try to reduce the use of this kind of instructions

```
1  void main()
2  {
3    // pre-processing part ...
4
5    // solver or computational part
6    // there is an implicit branch instruction in the loop
7    for(int i = 0; i < n; i++) {
8      // compiler generate a branch instruction here
9      if(i >= 1 && i < n -1) {
10       // compiler generate an other branch instruction here
11       switch(i % 4) {
12         case 0: B[i] = A[i] * 0.3333f;
13         case 1: B[i] = A[i] + 1.3333f;
14         case 2: B[i] = A[i] - 0.7555f;
15         case 3: B[i] = A[i] * 1.1111f;
16         default: break;
17       }
18     }
19   }
20
21   // post-processing part..
22 }
```

Computational code with branch instructions

# Reducing the number of branch instructions

```
1  void main()
2  {
3    // pre-processing part
4    const int n = 40000002;
5    float *A, *B;
6    A = new float[n]; // in
7    B = new float[n]; // out
8    randomInit(A, n);
9    struct timeval t1, t2;
10
11   // solver or computational part
12   for(int i = 1; i < n -1; i += 4) {
13     B[i+0] = A[i+0] + 1.3333f;
14     B[i+1] = A[i+1] - 0.7555f;
15     B[i+2] = A[i+2] * 1.1111f;
16     B[i+3] = A[i+3] * 0.3333f;
17   }
18
19   // post-processing part
20   delete[] A;
21   delete[] B;
22 }
```

Computational code with less branch instructions

- In the previous code with have replaced the `switch` by a 4-unrolling loop
- And we have replaced the `if` by starting the loop with `1` and finishing it with `n -1`

## Division

- Here is the cost of the main operations:
    - `add`: 1 CPU cycle
    - `sub`: 1 CPU cycle
    - `mul`: 1 CPU cycle
    - `div`: $\approx$ 20 CPU cycles
- As we can see, a division is very expensive compared to a multiplication
- It is much better to compute the inverse number and multiply by it!
- Be careful, when we multiply by inverse we lose some precision in the calculation

# Division example

```
1   void main()
2   {
3     // pre-processing part
4     int n = 1000;
5     float *A, *B;
6     A = new float[n]; // in
7     B = new float[n]; // out
8     randomInit(A, n);
9
10    // solver or computational part
11    for(int i = 0; i < n; i++)
12      B[i] = A[i] / 3.f;
13
14    // post-processing part
15    delete[] A;
16    delete[] B;
17  }
```

Simple computational code with divisions

- Theoretical number of cycles: $n \times 20$

# Division example

```
1   void main()
2   {
3     // pre-processing part
4     int n = 1000;
5     float *A, *B;
6     A = new float[n]; // in
7     B = new float[n]; // out
8     randomInit(A, n);
9     float inv3 = 1.f / 3.f;
10
11    // solver or computational part
12    for(int i = 0; i < n; i++)
13      B[i] = A[i] * inv3;
14
15    // post-processing part
16    delete[] A;
17    delete[] B;
18  }
```

Simple computational code without divisions

- Theoretical number of cycles: $20 + (n \times 1)$

# Special functions

- Here is the cost of the main special functions:
    - `pow`: very expensive, the number of cycles depends on the input
    - `sqrt`: ≈ 30 CPU cycles
    - `rsqrt`: 1 CPU cycle
    - `cos`: very expensive, the number of cycles depends on the input
    - `sin`: very expensive, the number of cycles depends on the input
    - `tan`: very expensive, the number of cycles depends on the input
- `rsqrt` take 1 cycle!
    - This is very surprising
    - In fact there are hardware pre-compute tables in today CPUs
    - The CPU simply returns the nearest value we need
- `pow`, `sqrt`, `cos`, `sin` and `tan` are very expensive try to not use them in the solver part of the code
- If it is not possible, at least try to reduce the number of calls

## Memory accesses

- When the code is memory bound, we have to carefully consider data structures and data accesses
- Memory bandwidth is slow compared to CPU computational capacity
    - There are some mechanisms to reduce this memory lack: pre-fetching data
    - Remember, cache accesses are done per line of words (and not per words)
    - So, it is very interesting to work on stream data
    - Also, we have to reduce direct accesses in RAM and maximize accesses in cache
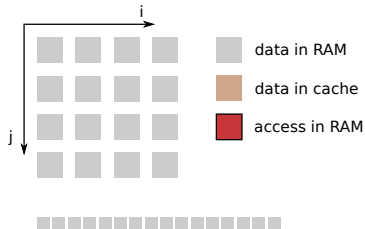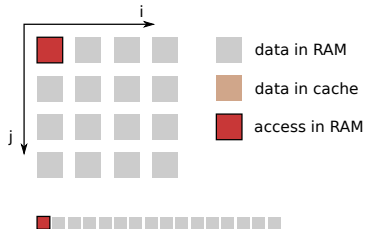
# Memory accesses example

```
1  void main()
2  {
3    // pre-processing part
4    int n = 4;
5    float *A, *B, *C;
6    A  = new float[n*n]; // in
7    B  = new float[n*n]; // in
8    C  = new float[n*n]; // out
9    randomInit(A, n); randomInit(B, n);
10
11   // solver or computational part
12   for(int i = 0; i < n; i++) // column
13     for(int j = 0; j < n; j++) // row
14       C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16   // post-processing part
17   delete[] A;
18   delete[] B;
19   delete[] C;
20 }
```

Adding square matrices



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory

# Memory accesses example

```
1  void main()
2  {
3    // pre-processing part
4    int n = 4;
5    float *A, *B, *C;
6    A  = new float[n*n]; // in
7    B  = new float[n*n]; // in
8    C  = new float[n*n]; // out
9    randomInit(A, n); randomInit(B, n);
10
11   // solver or computational part
12   for(int i = 0; i < n; i++) // column
13     for(int j = 0; j < n; j++) // row
14       C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16   // post-processing part
17   delete[] A;
18   delete[] B;
19   delete[] C;
20 }
```

Adding square matrices



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory
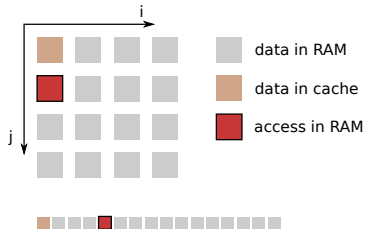
# Memory accesses example

```
 1  void main()
 2  {
 3    // pre-processing part
 4    int n = 4;
 5    float *A, *B, *C;
 6    A = new float[n*n]; // in
 7    B = new float[n*n]; // in
 8    C = new float[n*n]; // out
 9    randomInit(A, n); randomInit(B, n);
10
11    // solver or computational part
12    for(int i = 0; i < n; i++) // column
13      for(int j = 0; j < n; j++) // row
14        C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16    // post-processing part
17    delete[] A;
18    delete[] B;
19    delete[] C;
20  }
```

Adding square matrices



data in RAM

data in cache

access in RAM

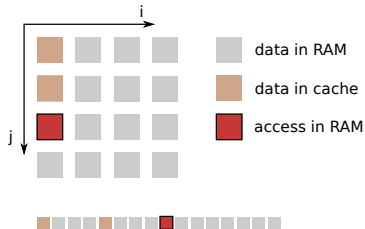Logical and hardware view of matrix in memory

# Memory accesses example

```
1  void main()
2  {
3    // pre-processing part
4    int n = 4;
5    float *A, *B, *C;
6    A  = new float[n*n]; // in
7    B  = new float[n*n]; // in
8    C  = new float[n*n]; // out
9    randomInit(A, n); randomInit(B, n);
10
11   // solver or computational part
12   for(int i = 0; i < n; i++) // column
13     for(int j = 0; j < n; j++) // row
14       C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16   // post-processing part
17   delete[] A;
18   delete[] B;
19   delete[] C;
20 }
```

Adding square matrices



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory
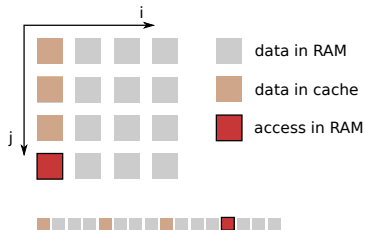
# Memory accesses example

```
1  void main()
2  {
3    // pre-processing part
4    int n = 4;
5    float *A, *B, *C;
6    A = new float[n*n]; // in
7    B = new float[n*n]; // in
8    C = new float[n*n]; // out
9    randomInit(A, n); randomInit(B, n);
10
11   // solver or computational part
12   for(int i = 0; i < n; i++) // column
13     for(int j = 0; j < n; j++) // row
14       C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16   // post-processing part
17   delete[] A;
18   delete[] B;
19   delete[] C;
20 }
```

Adding square matrices



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory
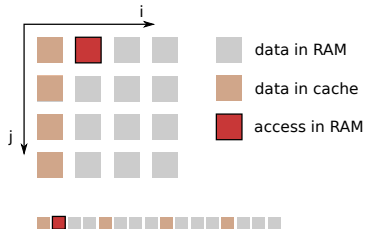
# Memory accesses example

```
 1  void main()
 2  {
 3    // pre-processing part
 4    int n = 4;
 5    float *A, *B, *C;
 6    A  = new float[n*n]; // in
 7    B  = new float[n*n]; // in
 8    C  = new float[n*n]; // out
 9    randomInit(A, n); randomInit(B, n);
10
11    // solver or computational part
12    for(int i = 0; i < n; i++) // column
13      for(int j = 0; j < n; j++) // row
14        C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16    // post-processing part
17    delete[] A;
18    delete[] B;
19    delete[] C;
20  }
```

Adding square matrices



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory
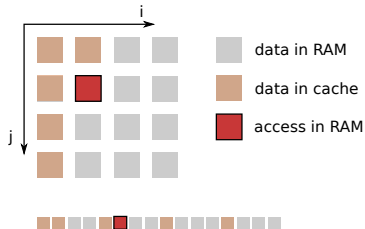
# Memory accesses example

```
 1  void main()
 2  {
 3    // pre-processing part
 4    int n = 4;
 5    float *A, *B, *C;
 6    A  = new float[n*n]; // in
 7    B  = new float[n*n]; // in
 8    C  = new float[n*n]; // out
 9    randomInit(A, n); randomInit(B, n);
10
11    // solver or computational part
12    for(int i = 0; i < n; i++) // column
13      for(int j = 0; j < n; j++) // row
14        C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16    // post-processing part
17    delete[] A;
18    delete[] B;
19    delete[] C;
20  }
```
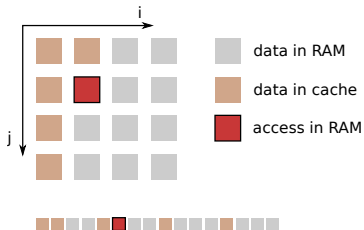
Adding square matrices



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory

# Memory accesses example



Logical and hardware view of matrix in memory

- In this implementation data accesses are not contiguous in memory
- There is a 4-stride between each access
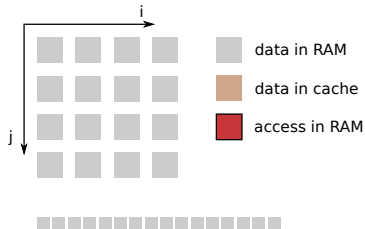
# Memory accesses example: solution

```
1  void main()
2  {
3    // pre-processing part
4    int n = 4;
5    float *A, *B, *C;
6    A  = new float[n*n]; // in
7    B  = new float[n*n]; // in
8    C  = new float[n*n]; // out
9    randomInit(A, n); randomInit(B, n);
10
11   // solver or computational part
12   for(int j = 0; j < n; j++) // row
13     for(int i = 0; i < n; i++) // column
14       C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16   // post-processing part
17   delete[] A;
18   delete[] B;
19   delete[] C;
20 }
```

Adding square matrices (solution)



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory
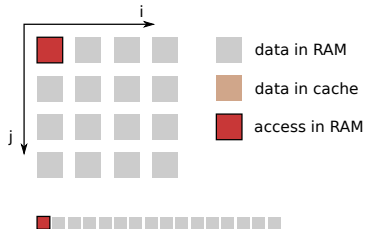
## Memory accesses example: solution

```
1   void main()
2   {
3     // pre-processing part
4     int n = 4;
5     float *A, *B, *C;
6     A  = new float[n*n]; // in
7     B  = new float[n*n]; // in
8     C  = new float[n*n]; // out
9     randomInit(A, n); randomInit(B, n);
10
11    // solver or computational part
12    for(int j = 0; j < n; j++) // row
13      for(int i = 0; i < n; i++) // column
14        C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16    // post-processing part
17    delete[] A;
18    delete[] B;
19    delete[] C;
20  }
```

Adding square matrices (solution)



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory
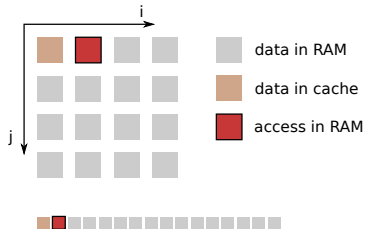
# Memory accesses example: solution

```
 1  void main()
 2  {
 3    // pre-processing part
 4    int n = 4;
 5    float *A, *B, *C;
 6    A  = new float[n*n]; // in
 7    B  = new float[n*n]; // in
 8    C  = new float[n*n]; // out
 9    randomInit(A, n); randomInit(B, n);
10
11    // solver or computational part
12    for(int j = 0; j < n; j++) // row
13      for(int i = 0; i < n; i++) // column
14        C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16    // post-processing part
17    delete[] A;
18    delete[] B;
19    delete[] C;
20  }
```

Adding square matrices (solution)



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory
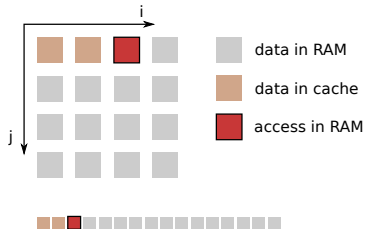
## Memory accesses example: solution

```
 1  void main()
 2  {
 3    // pre-processing part
 4    int n = 4;
 5    float *A, *B, *C;
 6    A  = new float[n*n]; // in
 7    B  = new float[n*n]; // in
 8    C  = new float[n*n]; // out
 9    randomInit(A, n); randomInit(B, n);
10
11    // solver or computational part
12    for(int j = 0; j < n; j++) // row
13      for(int i = 0; i < n; i++) // column
14        C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16    // post-processing part
17    delete[] A;
18    delete[] B;
19    delete[] C;
20  }
```

Adding square matrices (solution)



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory
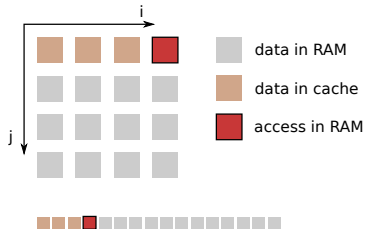
## Memory accesses example: solution

```
 1  void main()
 2  {
 3    // pre-processing part
 4    int n = 4;
 5    float *A, *B, *C;
 6    A = new float[n*n]; // in
 7    B = new float[n*n]; // in
 8    C = new float[n*n]; // out
 9    randomInit(A, n); randomInit(B, n);
10
11    // solver or computational part
12    for(int j = 0; j < n; j++) // row
13      for(int i = 0; i < n; i++) // column
14        C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16    // post-processing part
17    delete[] A;
18    delete[] B;
19    delete[] C;
20  }
```

Adding square matrices (solution)



- data in RAM
- data in cache
- access in RAM

Logical and hardware view of matrix in memory
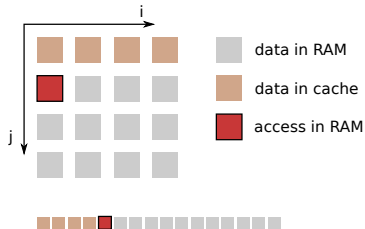
## Memory accesses example: solution

```
 1   void main()
 2   {
 3     // pre-processing part
 4     int n = 4;
 5     float *A, *B, *C;
 6     A  = new float[n*n]; // in
 7     B  = new float[n*n]; // in
 8     C  = new float[n*n]; // out
 9     randomInit(A, n); randomInit(B, n);
10
11     // solver or computational part
12     for(int j = 0; j < n; j++) // row
13       for(int i = 0; i < n; i++) // column
14         C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16     // post-processing part
17     delete[] A;
18     delete[] B;
19     delete[] C;
20   }
```

Adding square matrices (solution)



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory
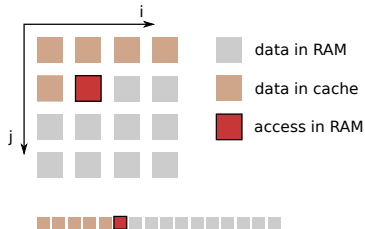
## Memory accesses example: solution

```
1   void main()
2   {
3     // pre-processing part
4     int n = 4;
5     float *A, *B, *C;
6     A  = new float[n*n]; // in
7     B  = new float[n*n]; // in
8     C  = new float[n*n]; // out
9     randomInit(A, n); randomInit(B, n);
10
11    // solver or computational part
12    for(int j = 0; j < n; j++) // row
13      for(int i = 0; i < n; i++) // column
14        C[i + j*n] = A[i + j*n] + B[i + j*n];
15
16    // post-processing part
17    delete[] A;
18    delete[] B;
19    delete[] C;
20  }
```
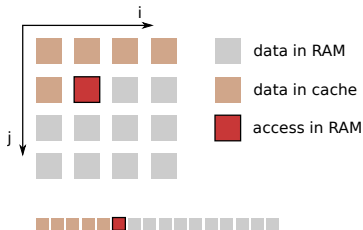
Adding square matrices (solution)



data in RAM

data in cache

access in RAM

Logical and hardware view of matrix in memory

## Memory accesses example: solution



Logical and hardware matrix memory view

- In this implementation data accesses are perfectly contiguous in memory
    - Loaded cache lines are well used
    - Performance will be better
- We have simply switched the i-loop with the j-loop
- In this case (addition of matrices) there is no data reuse
    - This code is very very memory bound!

# Cache blocking technique

- In many cases, some data can be reused!
- Let's take an example with a code working on a 2D grid

```
1   void main()
2   {
3     // pre-processing part
4     int cols = 10, rows = 6;
5     float *A = new float[cols*rows]; // in
6     float *B = new float[cols*rows]; // out
7     randomInit(A, cols*rows);
8
9     // solver or computational part
10    for(int j = 1; j < rows -1; j++) // row
11      for(int i = 1; i < cols -1; i++) // column
12        B[i + j*cols] = A[(i -1) + (j   )*cols] + A[(i +1) + (j   )*cols] +
13                        A[(i   ) + (j   )*cols] +
14                        A[(i   ) + (j -1)*cols] + A[(i   ) + (j +1)*cols];
15
16    // post-processing part
17    delete[] A;
18    delete[] B;
19  }
```
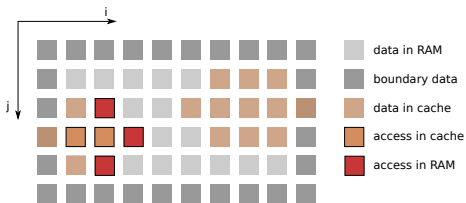
Stencil code

# Cache blocking technique

```
1  void main()
2  {
3    // pre-processing part ...
4
5    // solver or computational part
6    for(int j = 1; j < rows -1; j++) // row
7      for(int i = 1; i < cols -1; i++) // column
8        B[i + j*cols] = A[(i -1) + (j   )*cols] + A[(i +1) + (j   )*cols] + // left, right
9                        A[(i   ) + (j   )*cols] +                            // center
10                       A[(i   ) + (j -1)*cols] + A[(i   ) + (j +1)*cols];  // top, bottom
11
12   // post-processing part ...
13 }
```

Stencil code



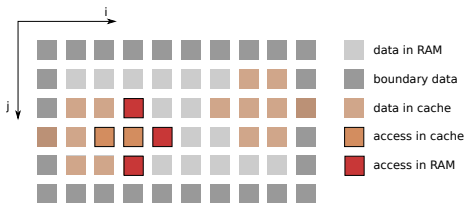| | |
|---|---|
| | data in RAM |
| | boundary data |
| | data in cache |
| | access in cache |
| | access in RAM |

Logical 2D grid memory view

# Cache blocking technique

```
1   void main()
2   {
3     // pre-processing part ...
4
5     // solver or computational part
6     for(int j = 1; j < rows -1; j++) // row
7       for(int i = 1; i < cols -1; i++) // column
8         B[i + j*cols] = A[(i -1) + (j   )*cols] + A[(i +1) + (j   )*cols] + // left, right
9                         A[(i   ) + (j   )*cols] +                           // center
10                        A[(i   ) + (j -1)*cols] + A[(i   ) + (j +1)*cols];  // top, bottom
11
12    // post-processing part ...
13  }
```

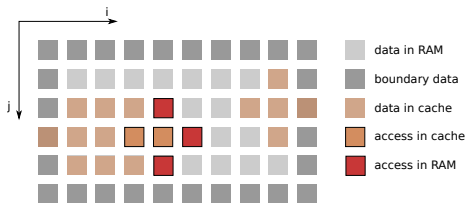Stencil code



Logical 2D grid memory view

# Cache blocking technique

```
1  void main()
2  {
3    // pre-processing part ...
4
5    // solver or computational part
6    for(int j = 1; j < rows -1; j++) // row
7      for(int i = 1; i < cols -1; i++) // column
8        B[i + j*cols] = A[(i -1) + (j   )*cols] + A[(i +1) + (j   )*cols] + // left, right
9                        A[(i   ) + (j   )*cols] +                           // center
10                       A[(i   ) + (j -1)*cols] + A[(i   ) + (j +1)*cols]; // top, bottom
11
12   // post-processing part ...
13 }
```

Stencil code



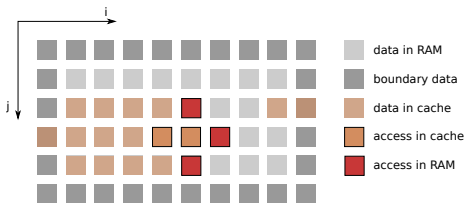| | data in RAM |
| | boundary data |
| | data in cache |
| | access in cache |
| | access in RAM |

Logical 2D grid memory view

# Cache blocking technique

```
1   void main()
2   {
3     // pre-processing part ...
4
5     // solver or computational part
6     for(int j = 1; j < rows -1; j++) // row
7       for(int i = 1; i < cols -1; i++) // column
8         B[i + j*cols] = A[(i -1) + (j  )*cols] + A[(i +1) + (j  )*cols] + // left, right
9                         A[(i  ) + (j  )*cols] +                           // center
10                        A[(i  ) + (j -1)*cols] + A[(i  ) + (j +1)*cols]; // top, bottom
11
12    // post-processing part ...
13  }
```
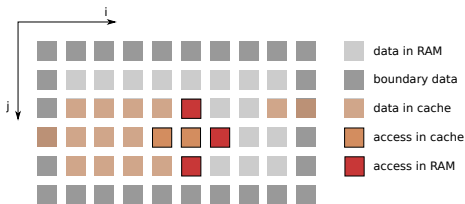
Stencil code



| | data in RAM |
| | boundary data |
| | data in cache |
| | access in cache |
| | access in RAM |

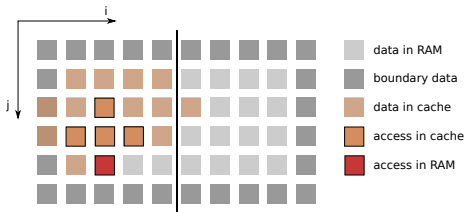Logical 2D grid memory view

## Cache blocking technique



Logical 2D grid memory view

- Each time we have 3 accesses in the RAM and 2 accesses in the cache
- Can we do better? Can we decrease the number of RAM accesses?
    - Yes, we can with the cache blocking technique!
    - The idea is to modify the data accessing manner in order to maximize data re-utilization

# Cache blocking technique: example



Logical 2D grid memory view

- With the cache blocking technique we reduce the number of RAM accesses
  - It remains just 1 access in RAM (sometimes 2)!
  - We cut the grid in different blocks (with a vertical separation here)

# Cache blocking technique: example



Logical 2D grid memory view

- With the cache blocking technique we reduce the number of RAM accesses
  - It remains just 1 access in RAM (sometimes 2)!
  - We cut the grid in different blocks (with a vertical separation here)
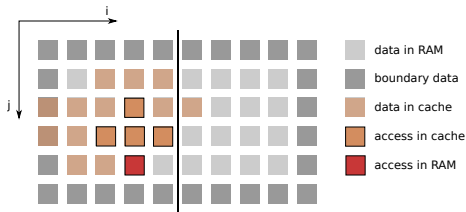
## Cache blocking technique: example



Logical 2D grid memory view

- With the cache blocking technique we reduce the number of RAM accesses
  - It remains just 1 access in RAM (sometimes 2)!
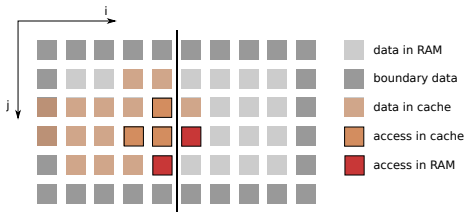  - We cut the grid in different blocks (with a vertical separation here)

# Cache blocking technique: example



Logical 2D grid memory view

- With the cache blocking technique we reduce the number of RAM accesses
  - It remains just 1 access in RAM (sometimes 2)!
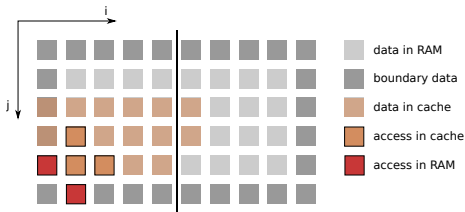  - We cut the grid in different blocks (with a vertical separation here)
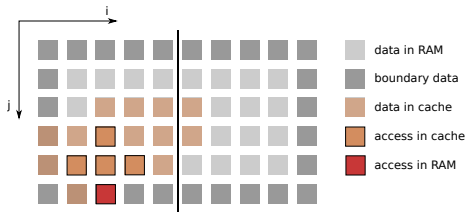
# Cache blocking technique: example



Logical 2D grid memory view

- With the cache blocking technique we reduce the number of RAM accesses
    - It remains just 1 access in RAM (sometimes 2)!
    - We cut the grid in different blocks (with a vertical separation here)

## Cache blocking technique

- How to define the size of blocks?
    - It depends on the problem
    - In previous stencil code the optimal block size can be computed like this:

$$blockSize = \frac{sizeOfCache}{2 \times 3 \times nThreads \times sizeOfData},$$

with $sizeOfCache$ the size of the biggest cache (L3) in bytes, $nThreads$ the number of threads we are using during the code execution and $sizeOfData$ the size of data we are computing (simple precision = 4 bytes, double precision = 8 bytes).

- We divide by 2 because the caches are optimal when we use half of them
- We divide by 3 because we have to keep 3 rows in cache with the stencil problem
- Be careful, if $blockSize \geq cols$ then cache blocking is useless

# Cache blocking technique: implementation

```
1   #define SIZE_OF_CACHE_L3 96 // bytes
2
3   void main()
4   {
5     // pre-processing part
6     int cols = 10, rows = 6;
7     float *A = new float[cols*rows]; // in
8     float *B = new float[cols*rows]; // out
9     randomInit(A, cols*rows);
10    int blockSize = SIZE_OF_CACHE_L3 / (2 * 3 * 1 * sizeof(float)); // (96 / 24) = 4
11
12    // solver or computational part
13    for(int iOff = 1; iOff < cols -1; iOff += blockSize) { // vertical blocks loop
14      blockSize = min(cols -1 - iOff, blockSize); // reduce the block size if necessary
15      for(int j = 1; j < rows -1; j++) // row
16        for(int i = iOff; i < iOff + blockSize; i++) // column
17          B[i + j*cols] = A[(i -1) + (j   )*cols] + A[(i +1) + (j   )*cols] +
18                          A[(i   ) + (j   )*cols] +
19                          A[(i   ) + (j -1)*cols] + A[(i   ) + (j +1)*cols];
20    }
21
22    // post-processing part
23    delete[] A;
24    delete[] B;
25  }
```

Stencil code with cache blocking

# Function calls

- A function call has a cost (extra assembly code)
- Is this is a sufficient reason to do not use functions a code?
    - It depends but sometimes yes it is!
    - In fact it depends on how many times we repeat the function call

```
1   void stencil(const float *A, float *B, const int i, const int j, const int cols)
2   {
3     B[i + j*cols] = A[(i -1) + (j   )*cols] + A[(i +1) + (j   )*cols] +
4                     A[(i   ) + (j   )*cols] +
5                     A[(i   ) + (j -1)*cols] + A[(i   ) + (j +1)*cols];
6   }
7
8   void main()
9   {
10    // pre-processing part
11    int cols = 10, rows = 6;
12    float *A = new float[cols*rows]; // in
13    float *B = new float[cols*rows]; // out
14    randomInit(A, cols*rows);
15
16    // solver or computational part
17    for(int j = 1; j < rows -1; j++) // row
18      for(int i = 1; i < cols -1; i++) // column
19        stencil(A, B, i, j, cols); // we call the stencil function many times!
20
21    // post-processing part
22    delete[] A;
23    delete[] B;
24  }
```

Stencil code with function calls

# What is inlining?

- Inlining a function is the same as replacing the function call by the code of the function itself
    - This way there is no more overhead because there are no more function calls
- We can manually do that but this is not a good idea...

# How to perform inlining?

- In term of software engineering, functions (or methods in object-oriented programming) are well spread
  - It is much better to use functions for the code readability
- So, can we build a beautiful and efficient code?
  - This is language or compiler dependent...
  - In C++ we have the `inline` keyword to perform inlining
  - In Fortran 90 (`ifort`) there is a directive: `!DEC$ ATTRIBUTES FORCEINLINE`
  - Often the compiler is free to perform inlining itself

```
1  inline void stencil(const float *A, float *B, const int i, const int j, const int cols) { ... }
2
3  void main()
4  {
5    // pre-processing part ...
6
7    // solver or computational part
8    for(int j = 1; j < rows -1; j++) // row
9      for(int i = 1; i < cols -1; i++) // column
10       stencil(A, B, i, j, cols); // we call the stencil function many times
11                                  // but with the inline keyword the compiler will
12                                  // automatically replace the call by the inner code
13
14   // post-processing part ...
15 }
```

Stencil code with function inlined calls in C++

## Work with the compiler

- Today compilers provide a lot of options to auto apply optimizations or to improve the performance of codes
- In this lesson we will talk about the C/C++ GNU compiler (gcc, g++) but you will find equivalent options with other compilers like the Intel compiler
- It is very important to know and understand what can and cannot be done by the compiler!
    - This way we can write beautiful (alias readable) and efficient codes
    - And the compiler can perform dirty optimizations when it generates the assembly code!

## Optimize options with GNU compiler

- The most famous option `-O[level]`:
  - `-O0`: reduces compilation time and makes debugging produce the expected results, this is the default.
  - `-O1`: the compiler tries to reduce code size and execution time, without performing any optimizations that take a great deal of compilation time.
  - `-O2`: optimizes even more, GCC performs nearly all supported optimizations that do not involve a space-speed trade-off.
  - `-O3`: optimizes even more, turns on the `-finline-functions`, `-funswitch-loops`, `-fpredictive-commoning`, `-fgcse-after-reload`, `-ftree-loop-vectorize`, `-ftree-loop-distribute-patterns`, `-ftree-slp-vectorize`, `-fvect-cost-model`, `-ftree-partial-pre` and `-fipa-cp-clone` options.
  - `-Ofast`: disregard strict standards compliance. It enables optimizations that are not valid for all standard-compliant programs. It turns on `-ffast-math` and the Fortran-specific `-fno-protect-parens` and `-fstack-arrays`.

# Specific optimize options with GNU compiler

- Most of the time we will avoid to use specific options
- But it is important to understand what can be performed by some of them:
    - `-finline-functions`: activate automatic inlining, the compiler is free to perform or not the optimization.
    - `-ftree-vectorize`: activate auto-vectorization of the code.
    - `-ffast-math`: do not respect IEEE specifications for the calculations (we lose some precision) but it can severely improve performances.
    - `-funroll-loops`: unroll loops whose number of iterations can be determined at compile time or upon entry to the loop. This option makes code larger, and may or may not make it run faster.
    - `-march=native`: allow specific instructions for a specific architecture, most of the time we will use this option in order to apply adapted vectorization on the code.
- The documentation of GNU optimize options: `https://gcc.gnu.org/onlinedocs/gcc/Optimize-Options.html`

# Contents

# Break instructions dependences

- Today CPUs need some independences between instructions
    - To fully use pipeline mechanism
    - And to efficiently exploit instruction-level parallelism (ILP)

```
1  void kernel(float *A, float *B, float *C, float *D, const float alpha, const int n)
2  {
3    for(int i = 0; i < n; i++) {
4      C[i] = A[i] + B[i];  // no dependences
5      D[i] = C[i] * alpha; // D depends on C
6      A[i] = D[i] - C[i];  // A depends on C and D
7      B[i] = A[i] * 2.0f;  // B depends on A
8    }
9  }
```

Kernel with a lot of dependences

# Break instructions dependences

```
1   void kernel(float *A, float *B, float *C, float *D, const float alpha, const int n)
2   {
3     for(int i = 0; i < n; i += 2) {
4       C[i+0] = A[i+0] + B[i+0]; // no dependences
5       C[i+1] = A[i+1] + B[i+1]; // no dependences
6
7       D[i+0] = C[i+0] * alpha;  // D[i+0] depends on C[i+0]
8       D[i+1] = C[i+1] * alpha;  // D[i+1] depends on C[i+1]
9
10      A[i+0] = D[i+0] - C[i+0]; // A[i+0] depends on C[i+0] and D[i+0]
11      A[i+1] = D[i+1] - C[i+1]; // A[i+1] depends on C[i+1] and D[i+1]
12
13      B[i+0] = A[i+0] * 2.0f;   // B[i+0] depends on A[i+0]
14      B[i+1] = A[i+1] * 2.0f;   // B[i+1] depends on A[i+1]
15    }
16  }
```

Kernel with less dependences

- Here, in the kernel example, it is possible to execute two instructions at the same time
- Sometimes compilers are capable of doing that on their own
- But when the code becomes complex the compiler can miss this kind of things

# Break instructions dependences

```c
#define UNROLL 2

void kernel(float *A, float *B, float *C, float *D, const float alpha, const int n)
{
  for(int i = 0; i < n; i += UNROLL) {
    for(int j = 0; j < UNROLL; j++)
      C[i+j] = A[i+j] + B[i+j]; // no dependences
    for(int j = 0; j < UNROLL; j++)
      D[i+j] = C[i+j] * alpha;  // D depends on C
    for(int j = 0; j < UNROLL; j++)
      A[i+j] = D[i+j] - C[i+j]; // A depends on C and D
    for(int j = 0; j < UNROLL; j++)
      B[i+j] = A[i+j] * 2.0f;   // B depends on A
  }
}
```

Kernel with less dependences ans less repetitions

- Previous code was quite ugly
- With the new one we can use the `-funroll-loops` compiler option (`gcc`) in order to unroll inner loops automatically
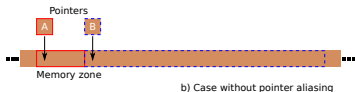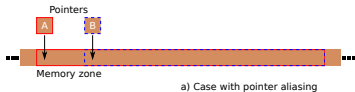- And there are no more code repetitions!

# Code vectorization

- Today vector instructions provide an efficient way to improve code performances
- How can we use those instructions?
    - By enabling the auto-vectorization process: the compiler automatically detects the zones to vectorize and generates an assembly code with vector instructions (be sure to have `-ftree-vectorize` option activated)
    - By calling specific functions: intrinsics
    - By directly writing assembly code: we will try to avoid this solution!
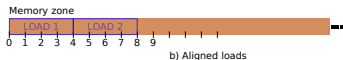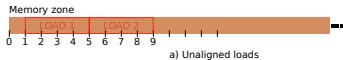
## Auto-vectorization

- Compilers are more and more capable of vectorizing codes automatically
- Loops are the best candidates for auto-vectorization
- This process requires respecting some constraints, the main idea is to write simple kernel codes (simple for the compiler)
  - The loop sizes have to be countable (`for`-loops without any `break`)
  - In general, branch statements (`if`, `switch`) are an obstacle to the vectorization
  - We have to guarantee that there is no pointer aliasing with the `__restrict` qualifier (this problem does not exist in `Fortran`)
    - There is pointer aliasing when two or more pointers can access the same memory zone
  - To achieve maximal bandwidth: loads and stores have to be aligned on the vector size

# Pointer aliasing and aligned loads/stores



Aliasing problem illustration



Alignment problem illustration

- a) A can access full B zone and B can access a sub-part of A
- b) A and B cannot cross

- a) Load 1 and 2 are not aligned on a multiple of 4
- b) Load 1 and 2 are well aligned on a multiple of 4

# How to verify if the code has been well vectorized?

- Most simple way is to take a look at the restitution time
- But sometimes we cannot assert things with the time alone
- GNU compilers provide vectorization reporting with the following option: `-ftree-vectorizer-verbose=[level]`
    - `0`: tells the compiler to report no diagnostic information.
    - `1`: tells the compiler to report on vectorized loops.
    - `2`: tells the compiler to report on vectorized and non-vectorized loops.
    - `3`: tells the compiler to report on vectorized and non-vectorized loops and any proven or assumed data dependences.
    - `4`: tells the compiler to report on non-vectorized loops.
    - `5`: tells the compiler to report on non-vectorized loops and the reason why they were not vectorized.

## Intrinsic calls

- Sometimes the compiler does not succeed in automatically vectorizing the code
- In this case we have to be more explicit and use intrinsic functions
- Basically an intrinsic call is equivalent to an assembly instruction
- This type of functions are very very hardware dependent!
    - This is why we will try not to use them unless we do not have the choice
- x86 intrinsics documentation: `https://software.intel.com/sites/landingpage/IntrinsicsGuide/`

# Intrinsic calls: example

```
1  // __restrict qualifier specify the compiler
2  // that there is no aliasing
3  void addVectors(const float* __restrict A,
4                  const float* __restrict B,
5                        float* __restrict C,
6                  const int n)
7  {
8    for(int i = 0; i < n; i++)
9      C[i] = A[i] + B[i];
10 }
```

Simple `addVectors` implementation

# Intrinsic calls: example

```
1  // __restrict qualifier specify the compiler
2  // that there is no aliasing
3  void addVectors(const float* __restrict A,
4                  const float* __restrict B,
5                        float* __restrict C,
6                  const int n)
7  {
8    for(int i = 0; i < n; i++)
9      C[i] = A[i] + B[i];
10 }
```

Simple `addVectors` implementation

- As you can see the intrinsics version is much more complex than the traditional version
- We have to limit intrinsics utilization for code readability

```
1  // headers for intrinsic AVX functions
2  #include "immintrin.h"
3
4  void iAddVectors(const float* __restrict A,
5                   const float* __restrict B,
6                         float* __restrict C,
7                   const int n)
8  {
9    // with AVX-256 we can compute vector of
10   // size 8 in single precision
11   for(int i = 0; i < n; i += 8) {
12     // load memory into vector registers
13     __m256 rA = _mm256_load_ps(A +i);
14     __m256 rB = _mm256_load_ps(B +i);
15
16     // perform SIMD/vectorized addition
17     __m256 rC = _mm256_add_ps(rA, rB);
18
19     // store C vector register into memory
20     _mm256_store_ps(C +i, rC);
21   }
22 }
```

Intrinsics AVX-256 `iAddVectors` implementation

# Contents

# Multi-core codes

- Multi-core architecture is well spread in the High Performance Computing
- There are two main ways to use multi-core architecture
    - Create multiple processes with MPI as a standard (distributed memory model)
    - Or create multiple threads with OpenMP as a standard (shared memory model)
- In this lesson we will not speak about multiple processes model
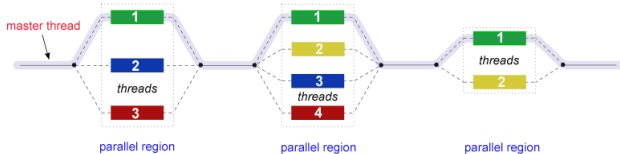- And we will go deeper into the multi-threaded model

# OpenMP

- OpenMP is a specific language for creating multi-threaded codes
- It is based on directives
  - Those directives describe how to perform the parallelism
  - The main advantage of directives is to not modify sequential code (in theory...)

```
 1  void addVectors(const float* __restrict A,
 2                  const float* __restrict B,
 3                        float* __restrict C,
 4                  const int n)
 5  {
 6  #pragma omp parallel // creation of a parallel zone directive (threads creation)
 7  {
 8  #pragma omp for // for-loop indices distribution directive
 9    for(int i = 0; i < n; i++)
10      C[i] = A[i] + B[i];
11  }
12  }
```
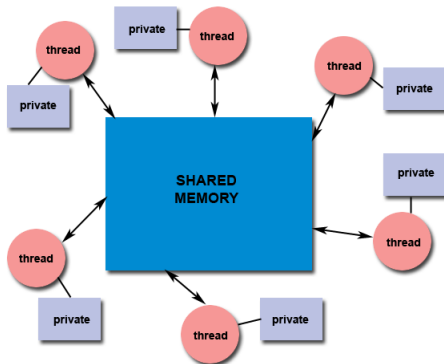
Simple `addVectors` OpenMP implementation

# OpenMP: fork-join model



Fork-join model illustration (Wikipedia)

- OpenMP follows the fork-join model
  - Each time we create a parallel zone (`#pragma omp parallel`) we create threads (fork operation)
  - At the end of a parallel zone threads are destroyed and there is an implicit barrier (join operation)
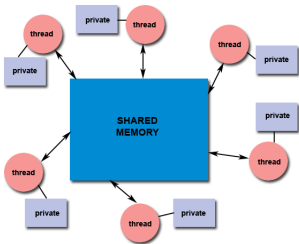    - Of course master thread remains

# OpenMP: shared memory model



Shared memory model illustration (computing.llnl.gov)

- OpenMP also follows the shared memory model
  - Each thread can access a global memory zone: the shared memory
  - But threads own also private data (not completely shared model)

# OpenMP: shared memory model example



Shared memory model illustration

```
1  void addVectors(const float* __restrict A,
2                  const float* __restrict B,
3                        float* __restrict C,
4                  const int n)
5  {
6  #pragma omp parallel
7  {
8  #pragma omp for
9    // i is private because it is declared
10   // after the omp parallel directive
11   for(int i = 0; i < n; i++)
12     // A, B and C are shared!
13     C[i] = A[i] + B[i];
14 }
15 }
```

Simple `addVectors` OpenMP implementation

# OpenMP: control data range

- OpenMP provides data range control
  - `private`: local to the thread,
  - `firstprivate`: local to the thread and initialized
  - `shared`: shared by all the threads, in `C/C++` this is the default
- Here `alpha` is a constant, we can put it in the private memory of each thread
- An efficient parallelism comes with minimal synchronisations
  - Shared data can generate a lot of synchronisations
  - Privacy increases thread independence

```
1   void dot(const float* __restrict A,
2            float* __restrict B,
3        const float alpha,
4        const int n)
5   {
6   #pragma omp parallel    \
7            shared(A, B) \
8            firstprivate(alpha, n)
9   {
10  #pragma omp for
11    for(int i = 0; i < n; i++)
12      B[i] = alpha * A[i];
13  }
14  }
```

OMP data range example

# OpenMP: for-loop indices distribution

- For-loop indices distribution can be controlled by the `schedule` clause
  - `static`: indices distribution is precomputed, and the amount of indices is the same for each thread
  - `dynamic`: indices distribution is done in real time along the loop execution, work load balancing can be better than with the `static` scheduling but `dynamic` scheduling costs some additional resources in order to attribute indices in real time
- In HPC, `static` distribution is the best choice if we are sure that each iteration has the same cost (in time)
- There are other types of scheduling but this is not a full OpenMP lesson

```
1  // ...
2  #pragma omp for schedule(static, 128) //we statistically attribute 128 per 128 indices to each threads
3    for(int i = 0; i < n; i++)
4      B[i] = alpha * A[i];
5  }
6  // ...
```

OMP scheduling example

# OpenMP: go further

- Previous slides were a brief overview of the main OpenMP principles
- To have more precise informations you can take a look at the very good OpenMP reference card:
  http://openmp.org/mp-documents/OpenMP-4.0-C.pdf
    - It could be a very good idea to print it and keep it ;-)
- In the next slides we will pay attention to OpenMP codes optimizations

# Avoid false sharing

- False sharing is a phenomena that occurs when threads write simultaneously data in a same line
  - Remember, the cache system works on lines of words: a line is the smallest element in caches coherence mechanism
    - If two or more threads are working on the same line they cannot write data simultaneously!
    - Stores are serialized and we talk about false sharing
- To avoid false sharing, threads have to work on a bigger amount of data than the cache line size
  - Concretely we have to avoid (static,1) or (dynamic,1) scheduling
  - Cache lines are not very big ($\approx$ 64 Bytes)
  - Just putting a (static,16) or (dynamic,16) resolves the problem
  - Be aware that default OpenMP scheduling is (static,1)!

## Reduce threads synchronisations: barriers

- In OpenMP there are a lot of implicit barriers, after each
    - `#pragma omp parallel` directive
    - `#pragma omp for` directive
    - `#pragma omp single` directive
- But not after `#pragma omp master` directive!
- If we are sure that there is no need to synchronise threads after the `#pragma omp for` directive, we can use the `nowait` clause
- Optimally we need only one `#pragma omp parallel` directive in a fully parallel code
    - OpenMP manages a pool of threads in order to reduce the cost of the `#pragma omp parallel` directive but this is not free, each time OpenMP has to reorganize the pool and wakes up the required threads

# Reduce threads synchronisations: barriers

```
 1  void kernelV1(const float *A, // size n
 2                const float *B, // size n
 3                const float *C, // size n
 4                      float *D, // size 2n
 5                const float alpha,
 6                const int n)
 7  {
 8  // threads creation overhead and
 9  // private variables creation overhead
10  #pragma omp parallel shared(A, B, D) \
11                        firstprivate(alpha, n)
12  {
13  #pragma omp for schedule(static,16)
14  {
15    for(int i = 0; i < n; i++)
16      D[i] = alpha * A[i] + B[i];
17  } // implicit barrier
18  } // implicit barrier
19
20  // threads attribution overhead and
21  // private variables creation overhead
22  #pragma omp parallel shared(A, C, D) \
23                        firstprivate(n)
24  {
25  #pragma omp for schedule(static,16)
26  {
27    for(int i = 0; i < n; i++)
28      D[n+i] = A[i] + C[i];
29  } // implicit barrier
30  } // implicit barrier
31  }
```

A lot of OMP barriers

# Reduce threads synchronisations: barriers

```
1  void kernelV1(const float *A, // size n
2                const float *B, // size n
3                const float *C, // size n
4                      float *D, // size 2n
5                const float alpha,
6                const int n)
7  {
8  // threads creation overhead and
9  // private variables creation overhead
10 #pragma omp parallel shared(A, B, D) \
11                      firstprivate(alpha, n)
12 {
13 #pragma omp for schedule(static,16)
14 {
15   for(int i = 0; i < n; i++)
16     D[i] = alpha * A[i] + B[i];
17 } // implicit barrier
18 } // implicit barrier
19
20 // threads attribution overhead and
21 // private variables creation overhead
22 #pragma omp parallel shared(A, C, D) \
23                      firstprivate(n)
24 {
25 #pragma omp for schedule(static,16)
26 {
27   for(int i = 0; i < n; i++)
28     D[n+i] = A[i] + C[i];
29 } // implicit barrier
30 } // implicit barrier
31 }
```

A lot of OMP barriers

```
1  void kernelV2(const float *A, // size n
2                const float *B, // size n
3                const float *C, // size n
4                      float *D, // size 2n
5                const float alpha,
6                const int n)
7  {
8  // threads creation overhead and
9  // private variables creation overhead
10 #pragma omp parallel shared(A, B, C, D) \
11                      firstprivate(alpha, n)
12 {
13 #pragma omp for schedule(static,16) nowait
14 {
15   for(int i = 0; i < n; i++)
16     D[i] = alpha * A[i] + B[i];
17 } // no implicit barrier (nowait clause)
18
19 #pragma omp for schedule(static,16)
20 {
21   for(int i = 0; i < n; i++)
22     D[n+i] = A[i] + C[i];
23 } // implicit barrier
24 } // implicit barrier
25 }
```

Less OMP barriers

# Reduce threads synchronisations: critical sections

- Sometimes it is not possible to have a fully parallel code and some regions of the code remain intrinsically sequential
- In OpenMP we can specify this kind of region with the `#pragma omp critical` directive
- But we have to use this directive carefully
  - It can be a main cause of slow down in OpenMP codes!

# Reduce threads synchronisations: critical sections

```
1  float kernelV1(const float *A, // size n
2                 const float *B, // size n
3                 const int n)
4  {
5    float minVal = INF;
6
7  #pragma omp parallel shared(A, B, minVal) \
8                       firstprivate(n)
9  {
10 #pragma omp for schedule(static,16)
11 {
12   for(int i = 0; i < n; i++) {
13     B[i] = 0.5f * A[i];
14
15 #pragma omp critical // we are sure that only
16                      // one thread can
17                      // modify minVal
18 {
19     if(B[i] < minVal)
20       minVal = B[i];
21 }
22   }
23 }
24 }
25
26   return minVal;
27 }
```

Critical section

# Reduce threads synchronisations: critical sections

```
1  float kernelV1(const float *A, // size n
2                 float *B, // size n
3                 const int n)
4  {
5    float minVal = INF;
6
7  #pragma omp parallel shared(A, B, minVal) \
8                      firstprivate(n)
9  {
10 #pragma omp for schedule(static,16)
11 {
12   for(int i = 0; i < n; i++) {
13     B[i] = 0.5f * A[i];
14
15 #pragma omp critical // we are sure that only
16                      // one thread can
17                      // modify minVal
18 {
19     if(B[i] < minVal)
20        minVal = B[i];
21 }
22   }
23 }
24 }
25
26   return minVal;
27 }
```

Critical section

```
1  float kernelV2(const float *A, // size n
2                 float *B, // size n
3                 const int n)
4  {
5    float minVal = INF;
6
7  #pragma omp parallel shared(A, B, minVal) \
8                      firstprivate(n)
9  {
10 #pragma omp for schedule(static,16)
11 {
12   for(int i = 0; i < n; i++) {
13     B[i] = 0.5f * A[i];
14
15     // when B[i] >= minVal there is no more
16     // threads synchronisation
17     if(B[i] < minVal)
18 #pragma omp critical
19 {
20       // this is very important to re-do the
21       // test because an other thread may
22       // have modify the minVal value
23       if(B[i] < minVal)
24         minVal = B[i];
25 }
26   }
27 }
28 }
29
30   return minVal;
31 }
```

Critical section optimized

# Search algorithms

- In OpenMP 3 there is no optimal solution for search algorithms
- This kind of algorithm typically requires while-loops or do-while-loops
- However there is a tip to fix this lack in OpenMP 3
- Latest version of OpenMP (v4) provides better control of threads
  - We can terminate threads...
  - We will not speak about OpenMP 4 because many current systems does not support this version

# Search algorithms: OpenMP 3 tip

```
1   bool searchValV1(const float *A,
2                    const int n,
3                        float val)
4   {
5     bool found = false;
6
7   #pragma omp parallel shared(A, found) \
8                       firstprivate(val)
9   {
10  #pragma omp for schedule(static,16)
11  {
12    for(int i = 0; i < n; i++) {
13      if(A[i] == val)
14        found = true;
15    }
16  }
17  }
18
19    return found;
20  }
```

Search algorithms

# Search algorithms: OpenMP 3 tip

```
1  bool searchValV1(const float *A,
2                   const int n,
3                   float val)
4  {
5    bool found = false;
6
7  #pragma omp parallel shared(A, found) \
8                   firstprivate(val)
9  {
10 #pragma omp for schedule(static,16)
11 {
12   for(int i = 0; i < n; i++) {
13     if(A[i] == val)
14       found = true;
15   }
16 }
17 }
18
19   return found;
20 }
```

Search algorithms

```
1  bool searchValV2(const float *A,
2                   const int n,
3                   float val)
4  {
5    bool found = false;
6
7  #pragma omp parallel shared(A, found) \
8                   firstprivate(val)
9  {
10 #pragma omp for schedule(static,16)
11 {
12   for(int i = 0; i < n; i++) {
13     if(!found) // we are doing nothing if
14               // we have found the value
15               // in the array
16       if(A[i] == val)
17         found = true;
18   }
19 }
20 }
21
22   return found;
23 }
```

Search algorithms optimized

## Final words

- Of course there are many more possible optimizations
- The purpose of this lesson was to raise awareness among optimization problematic
- Now you have tools to understand and to create your own optimizations
- Presented techniques are very often specific to the problem nature
    - Be aware that there is no perfect optimization
    - You will have to think about your needs before trying to optimize your code