

## Lecture 22

# Writing Cache-friendly Code

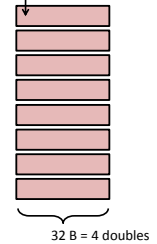
CPSC 275  
Introduction to Computer Systems

## A Higher Level Example

```
int sum_array_rows(double a[16][16])
{
    int i, j;
    double sum = 0;
    for (i = 0; i < 16; i++)
        for (j = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```

```
int sum_array_cols(double a[16][16])
{
    int i, j;
    double sum = 0;
    for (j = 0; j < 16; j++)
        for (i = 0; i < 16; i++)
            sum += a[i][j];
    return sum;
}
```

Ignore the variables sum, i, j  
assume: cold (empty) cache,  
a[0][0] goes here



## Writing Cache Friendly Code

- Make the common case go fast
  - Focus on the inner loops of the core functions
- Minimize the misses in the inner loops
  - Repeated references to variables are good (**temporal locality**)
  - Stride-1 reference patterns are good (**spatial locality**)

## Matrix Multiplication Example

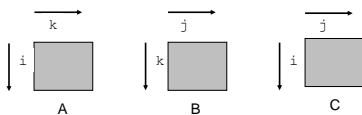
- Description:
  - Multiply N x N matrices
  - $O(N^3)$  total operations
  - N reads per source element
  - N values summed per destination
    - but may be able to hold in register

```
/* ijk */
for (i=0; i<n; i++) {
    for (j=0; j<n; j++) {
        sum = 0.0;
        for (k=0; k<n; k++)
            sum += a[i][k] * b[k][j];
        c[i][j] = sum;
    }
}
```

Variable sum  
held in register

## Miss Rate Analysis for Matrix Multiply

- Assume:
  - Line size = 32 bytes (big enough for 4 doubles)
  - Matrix dimension (N) is very large
  - Cache is not even big enough to hold multiple rows
- Analysis Method:
  - Look at access pattern of inner loop



## Layout of C Arrays in Memory

- C arrays allocated in *row-major* order
  - each row in contiguous memory locations
- Stepping through columns in one row:
 

```
for (i = 0; i < n; i++)
    sum += a[0][i];
```

  - accesses successive elements
  - exploits spatial locality
- Stepping through rows in one column:
 

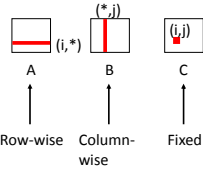
```
for (i = 0; i < n; i++)
    sum += a[i][0];
```

  - accesses distant elements
  - no spatial locality!

## Matrix Multiplication (ijk)

```
/* ijk */
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}
```

Inner loop:



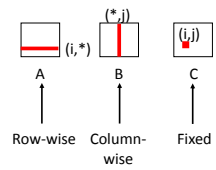
Misses per inner loop iteration:

A	B	C
0.25	1.0	0.0

## Matrix Multiplication (jik)

```
/* jik */
for (j=0; j<n; j++) {
  for (i=0; i<n; i++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}
```

Inner loop:



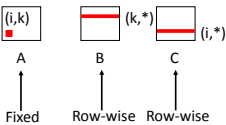
Misses per inner loop iteration:

A	B	C
0.25	1.0	0.0

## Matrix Multiplication (kij)

```
/* kij */
for (k=0; k<n; k++) {
  for (i=0; i<n; i++) {
    r = a[i][k];
    for (j=0; j<n; j++)
      c[i][j] += r * b[k][j];
  }
}
```

Inner loop:



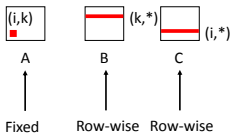
Misses per inner loop iteration:

A	B	C
0.0	0.25	0.25

## Matrix Multiplication (ikj)

```
/* ikj */
for (i=0; i<n; i++) {
  for (k=0; k<n; k++) {
    r = a[i][k];
    for (j=0; j<n; j++)
      c[i][j] += r * b[k][j];
  }
}
```

Inner loop:



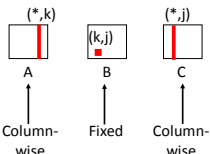
Misses per inner loop iteration:

A	B	C
0.0	0.25	0.25

## Matrix Multiplication (jki)

```
/* jki */
for (j=0; j<n; j++) {
  for (k=0; k<n; k++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
```

Inner loop:



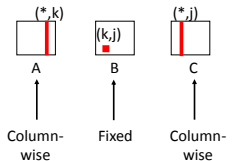
Misses per inner loop iteration:

A	B	C
1.0	0.0	1.0

## Matrix Multiplication (kji)

```
/* kji */
for (k=0; k<n; k++) {
  for (j=0; j<n; j++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
```

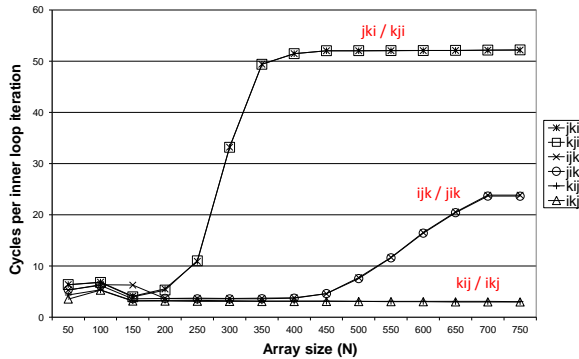
Inner loop:



Misses per inner loop iteration:

A	B	C
1.0	0.0	1.0

## Core i7 Matrix Multiply Performance

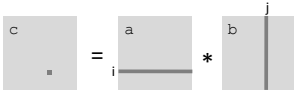


## Optimizations for Memory Hierarchy

- Write code that has locality
  - Spatial: access data contiguously
  - Temporal: make sure access to the same data is not too far apart in time
- How to achieve?
  - Proper choice of algorithm
  - Loop transformations

## Example: Matrix Multiplication

```
c = (double *) calloc(sizeof(double), n*n);
/* Multiply n x n matrices a and b */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i++)
        for (j = 0; j < n; j++)
            for (k = 0; k < n; k++)
                c[i*n + j] += a[i*n + k]*b[k*n + j];
}
```

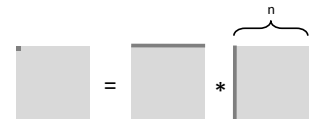


## Cache Miss Analysis

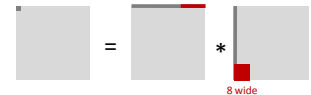
- Assume:
  - Matrix elements are doubles
  - Cache block = 64 bytes = 8 doubles
  - Cache size  $C \ll n$  (much smaller than  $n$ )

### First iteration:

- $n/8 + n = 9n/8$  misses (omitting matrix c)



- Afterwards **in cache**: (schematic)

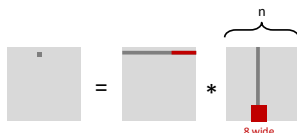


## Cache Miss Analysis

- Assume:
  - Matrix elements are doubles
  - Cache block = 64 bytes = 8 doubles
  - Cache size  $C \ll n$  (much smaller than  $n$ )

### Other iterations:

- Again:  $n/8 + n = 9n/8$  misses (omitting matrix c)

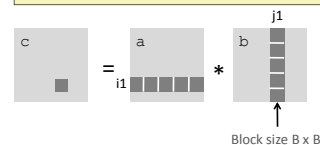


### Total misses:

- $9n/8 * n^2 = (9/8) * n^3$

## Blocked Matrix Multiplication

```
c = (double *) calloc(sizeof(double), n*n);
/* Multiply n x n matrices a and b */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i+=B)
        for (j = 0; j < n; j+=B)
            for (k = 0; k < n; k+=B)
                /* B x B mini matrix multiplications */
                for (il = i; il < i+B; il++)
                    for (jl = j; jl < j+B; jl++)
                        for (kl = k; kl < k+B; kl++)
                            c[il*n + jl] += a[il*n + kl]*b[kl*n + jl];
}
```



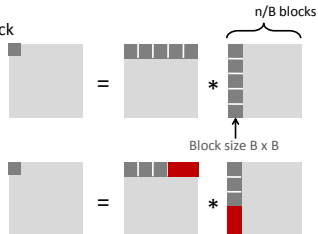
## Cache Miss Analysis

- Assume:
  - Cache block = 64 bytes = 8 doubles
  - Cache size  $C \ll n$  (much smaller than  $n$ )
  - Three blocks  $\blacksquare$  fit into cache:  $3B^2 < C$

- First (block) iteration:

- $B^2/8$  misses for each block
  - $2n/B * B^2/8 = nB/4$  (omitting matrix  $c$ )

- Afterwards in cache (schematic)



## Cache Miss Analysis

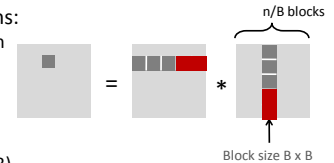
- Assume:
  - Cache block = 64 bytes = 8 doubles
  - Cache size  $C \ll n$  (much smaller than  $n$ )
  - Three blocks  $\blacksquare$  fit into cache:  $3B^2 < C$

- Other (block) iterations:

- Same as first iteration
  - $2n/B * B^2/8 = nB/4$

- Total misses:

- $nB/4 * (n/B)^2 = n^3/(4B)$



## Summary

- No blocking:  $(9/8) * n^3$
- Blocking:  $1/(4B) * n^3$
- If  $B = 8$  difference is  $4 * 8 * 9 / 8 = 36x$
- If  $B = 16$  difference is  $4 * 16 * 9 / 8 = 72x$
- Suggests largest possible block size  $B$ , but limit  $3B^2 < C$ !
- Reason for dramatic difference:
  - Matrix multiplication has inherent temporal locality:
    - Input data:  $3n^2$ , computation  $2n^3$
    - Every array element used  $O(n)$  times!
  - But program has to be written properly

## Cache-Friendly Code

- Programmer can optimize for cache performance
  - How data structures are organized
  - How data are accessed
    - Nested loop structure
    - Blocking is a general technique
- All systems favor "cache-friendly code"
  - Getting absolute optimum performance is very platform specific
    - Cache sizes, line sizes, associativities, etc.
  - Can get most of the advantage with generic code
    - Keep working set reasonably small (temporal locality)
    - Use small strides (spatial locality)
    - Focus on inner loop code

## Practice Problems

- Read CSaPP Sec. 6.4 and try 6.18 and 6.19