

Recap of the Last Class

- Cache organization and operation
 - Copy of data, usually in a smaller and faster location
 - Cache memory organization
 - Data managed in blocks (lines)
 - Memory address bits use to determine mapping
 - Direct-mapped, fully-associative, set-associative
 - Handling writes: write-back/write-through, write-allocate/write-no-allocate
 - Performance metrics: hits vs. misses, latency of access

35

Breakout

- What is s , b , t for an 8 Kbyte cache with 64-byte lines and a set-associativity of 4? Assume a 32-bit address space

36

Today

- Cache organization and operation
 - Quiz 6 out on blackboard
- Performance impact of caches
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

37

Breakout

- Assume a direct-mapped cache of size 8 KByte 64-Byte line. A is located at 0x8000 and C is located at 0x9000. How many misses would you incur? Would this change with a 2-way set associative cache? Assume i is retained in a register and the arrays are not initially in the cache.

```
int A[100], C[100];
for (i = 0; i < 100; i++)
    A[i] = C[i];
```

38

The Memory Mountain

- **Read throughput** (read bandwidth)
 - Number of bytes read from memory per second (MB/s)
- **Memory mountain:** Measured read throughput as a function of spatial and temporal locality
 - Compact way to characterize memory system performance

39

Memory Mountain Test Function

```

long data[MAXELEMS]; /* Global array to traverse */

/* test - Iterate over first "elems" elements of
 * array "data" with stride of "stride", using
 * using 4x4 loop unrolling.
 */
int test(int elems, int stride) {
    long i, sx2=stride*2, sx3=stride*3, sx4=stride*4;
    long acc0 = 0, acc1 = 0, acc2 = 0, acc3 = 0;
    long length = elems, limit = length - sx4;

    /* Combine 4 elements at a time */
    for (i = 0; i < limit; i += sx4) {
        acc0 = acc0 + data[i];
        acc1 = acc1 + data[i+stride];
        acc2 = acc2 + data[i+sx2];
        acc3 = acc3 + data[i+sx3];
    }

    /* Finish any remaining elements */
    for (; i < length; i++) {
        acc0 = acc0 + data[i];
    }

    return ((acc0 + acc1) + (acc2 + acc3));
}
    
```

mountain/mountain.c

Call test () with many combinations of elems and stride.

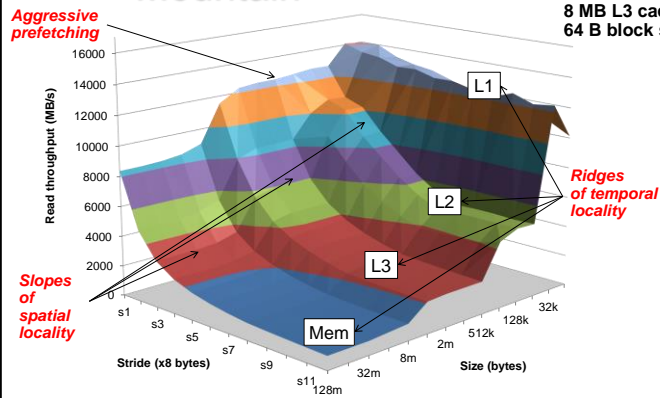
For each elems and stride:

1. Call test () once to warm up the caches.
2. Call test () again and measure the read throughput (MB/s)

40

The Memory Mountain

Core i7 Haswell
2.1 GHz
32 KB L1 d-cache
256 KB L2 cache
8 MB L3 cache
64 B block size



41

Today

- Cache organization and operation
- Performance impact of caches
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

42

Matrix Multiplication Example

- Description:
 - Multiply $N \times N$ matrices
 - Matrix elements are doubles (8 bytes)
 - $O(N^3)$ total operations
 - N reads per source element
 - N values summed per destination
 - but may be able to hold in register

```

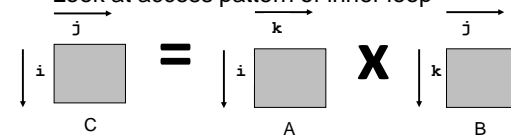
/* ijk */
for (i=0; i<n; i++) {
    for (j=0; j<n; j++) {
        sum = 0.0;
        for (k=0; k<n; k++)
            sum += a[i][k] * b[k][j];
        c[i][j] = sum;
    }
}
matmult/mm.c
    
```

Variable sum held in register

43

Miss Rate Analysis for Matrix Multiply

- Assume:
 - Block size = 32B (big enough for four doubles)
 - Matrix dimension (N) is very large
 - Approximate $1/N$ as 0.0
 - Cache is not even big enough to hold multiple rows
- Analysis Method:
 - Look at access pattern of inner loop



44

Layout of C Arrays in Memory (review)

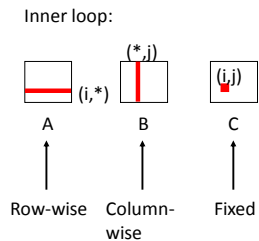
- C arrays allocated in row-major order
 - each row in contiguous memory locations
- Stepping through columns in one row:
 - for ($i = 0; i < N; i++$)
 - $sum += a[0][i];$
 - accesses successive elements
 - if block size (B) > $sizeof(a_{ij})$ bytes, exploit spatial locality
 - miss rate = $sizeof(a_{ij}) / B$
- Stepping through rows in one column:
 - for ($i = 0; i < n; i++$)
 - $sum += a[i][0];$
 - accesses distant elements
 - no spatial locality!
 - miss rate = 1 (i.e. 100%)

45

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;
    }
}
matmult/mm.c
    
```



Misses per inner loop iteration:

A	B	C
0.25	1.0	0.0

46

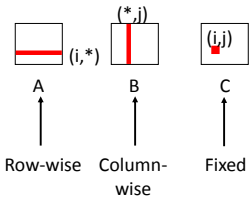
Matrix Multiplication (jik)

```

/* 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
  }
}
matmult/mm.c

```

Inner loop:



Misses per inner loop iteration:

A	B	C
0.25	1.0	0.0

47

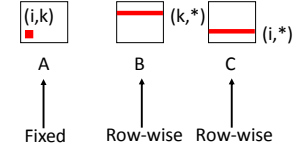
Matrix Multiplication (kij)

```

/* 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];
  }
}
matmult/mm.c

```

Inner loop:



Misses per inner loop iteration:

A	B	C
0.0	0.25	0.25

48

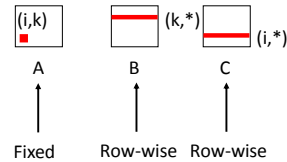
Matrix Multiplication (ikj)

```

/* 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];
  }
}
matmult/mm.c

```

Inner loop:



Misses per inner loop iteration:

A	B	C
0.0	0.25	0.25

49

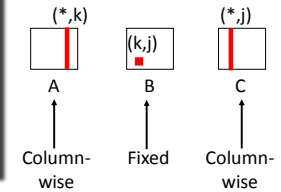
Matrix Multiplication (jki)

```

/* 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;
  }
}
matmult/mm.c

```

Inner loop:



Misses per inner loop iteration:

A	B	C
1.0	0.0	1.0

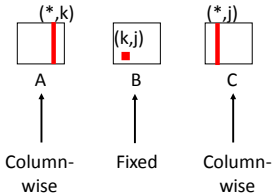
50

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;
  }
}
matmult/mm.c
    
```

Inner loop:



Misses per inner loop iteration:

A	B	C
1.0	0.0	1.0

51

Summary of Matrix Multiplication

```

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;
  }
}
    
```

ijk (& jik):

- 2 loads, 0 stores
- misses/iter = 1.25

```

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];
  }
}
    
```

kij (& ikj):

- 2 loads, 1 store
- misses/iter = 0.5

```

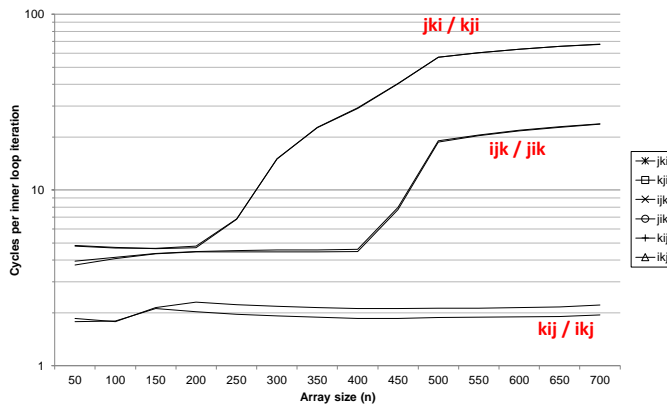
for (j=0; j<n; j++) {
  for (k=0; k<n; k++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
    
```

jki (& kji):

- 2 loads, 1 store
- misses/iter = 2.0

52

Core i7 Matrix Multiply Performance



53

Today

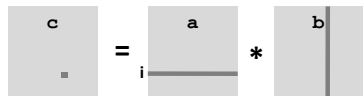
- Cache organization and operation
- Performance impact of caches
 - The memory mountain
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

54

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];
}
```



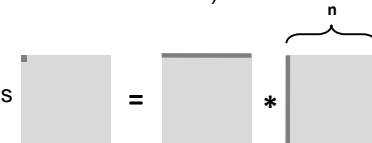
55

Cache Miss Analysis

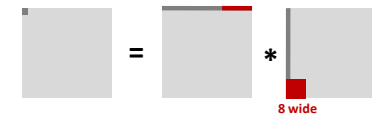
- Assume:
 - Matrix elements are doubles
 - Cache block = 8 doubles
 - Cache size $C \ll n$ (much smaller than n)

- First iteration:

- $n/8 + n = 9n/8$ misses



- Afterwards in cache: (schematic)



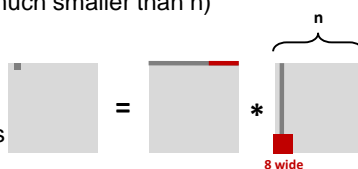
56

Cache Miss Analysis

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

- Second iteration:

- Again: $n/8 + n = 9n/8$ misses



- Total misses:

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

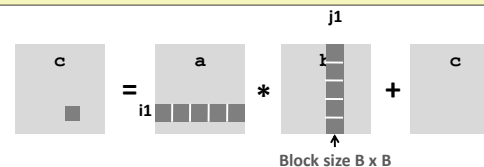
57

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 (i1 = i; i1 < i+B; i1++)
                    for (j1 = j; j1 < j+B; j1++)
                        for (k1 = k; k1 < k+B; k1++)
                            c[i1*n+j1] += a[i1*n+k1]*b[k1*n+j1];
}
```

matmult/bmm.c



58

Cache Miss Analysis

- Assume:
 - Cache block = 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)
-

59

Cache Miss Analysis

- Assume:
 - Cache block = 8 doubles
 - Cache size $C \ll n$ (much smaller than n)
 - Three blocks \blacksquare fit into cache: $3B^2 < C$
 - Second (block) iteration:
 - Same as first iteration
 - $2n/B * B^2/8 = nB/4$
 - Total misses:
 - $nB/4 * (n/B)^2 = n^3/(4B)$
-

60

Blocking Summary

- No blocking: $(9/8) * n^3$
- Blocking: $1/(4B) * n^3$
- Suggest 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 elements used $O(n)$ times!
 - But program has to be written properly

61

Concluding Observations

- 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)

63