Pro	gramming for performance on multi-core and many-core	22 <sup>nd</sup> March 2012	
Qu	z: 25 marks		[1 hour]
1.\	Vhy <b>Loop unrolling</b> is often useful in improving the performan	ce of numerical code?	[1m]
a) b) c)	It helps in having less instruction cache misses. It enables other optimizations, for example scalar replacement It enables better instruction scheduling.	nt or reassociation.	
Pic	k the true statements [ , , ].		
2.0	Consider the following code:	Ι	[2m]
dou for x[i] } The	ible * func (const double N, const double *x, int size) { (i=1, i < size; i++) = (x[i]/size) * N; e previous loop can be optimized as follows:		
dou dou for x[i] }	<pre>ible * func (const double N, const double *x, int size) { ible tmp = N/size; (i=1, i &lt; size; i++) = x[i] * tmp;</pre>		
DU			
<b>3.</b> A a)t b)t c)t d)t	GPU SM can load and store to and from (circleall that apply): he GPU's device memory (on-board, off-chip) he SM's local shared memory (on-chip) he shared memories of other SM's he L1 caches of other SMs	: []	[1m]
<b>4.</b> ( a) ( b) ( c) ( d) ( e) (	GPU thread divergence happens if (circle one): [ lifferent SMs execute different code different thread blocks in one SM execute different code lifferent warps in one thread block execute different code none of the above all of the above	]	[1m]

5. The code below was optimized in different ways resulting in versions v1 and v2. Both versions were tested on arrays of different sizes in the same computer. When v1 was used to compute on arrays of size 2000, the execution time was 64,000 cycles. When v2 was used when arrays of size 100, the execution time was 4000 cycles. Which version do you think is better optimized? Justify your answer. [2m]

```
for (int i = 0; i < size; i++){
C[i] = sqrt(A[i])+sqrt(B[i]);
}</pre>
```

Does a commercial compiler (icc) with higher optimization flags vectorize this loop? If not why? If so, how?
 [3m]

```
loc=maxint;
for (i=0; i<n; i++)
    if x[i]<0 {
        loc=i;
        break;
    }
```

Optimize the cache performance of the following code. Do not parallelize it. Do not consider prefetching nor TLB behavior. You can assume that none of the operations overflow. Assume that N and M are extremely large powers of 2. Assume that A[], B[], and C[] are integer arrays, and that an int is 4B.

```
for (i=0;i<N;i++)
for (j=0;j<M;j++)
A[i] += B[j][j] * C[i];</pre>
```

a) The following data structure Data2 is used to store one million samples of (x,y,z) 3D positioning data. How many bytes of memory will Data2 occupy?

```
struct mystruct {
    int x; // 4
    int y; // 4
    int z; // 4
} Data2[1000000];
```

b) How many bytes of memory will the following alternative data structure Data3 occupy?

```
struct mystruct {
    int x[1000000]; // int is 4
    int y[1000000]; // int is 4
    int z[1000000]; // int is 4
} Data3;
```

c) Data2 and Data3 above can both be used to store the same set of data, but would have different layouts in memory. If you were asked to write code to compute the distance from the origin of (x,y,z) for 1000 randomly-selected samples, which of Data2 or Data3 would be the preferred data structure and why is it better than the alternative?

d) If you were asked to write code to compute the mean of y across all million values, which of Data2 or Data3 would be the preferred data structure and why?

9. Generally, in which case would function inlining be least likely to harm performance: A) a large function called from many locations; B) a large function called from a few locations; C) a small function called from many locations; D) a small function called from a few locations. []

[3m]

```
11. Tiling: Given this machine:
1<sup>st</sup> level cache
128B cache blocks
8-way set associative
16 sets,
Page size: 8KB,
TLB: fully associative, 128 entries
```

If you were to tile the following code

// assume D[] is already initialized to zeros

Ignoring storage for instructions and small variables, and assuming that there are no unlucky conflict misses (only capacity misses), what is the largest tile size T (in number of elements, not bytes) that you can use to tile to minimize misses for the first level data cache?

## 12. Compiler Optimization

Name 8 optimizations that you would hope your compiler would do to this code. For each, give the formal name of the optimization and very briefly describe which variable(s) or code parts it will modify/improve and how. The first answer of 8 is given as an example.

```
1:x=5;
2:y=2;
3: debug = 0;
4: z=x+y;
5:
6: for (i=0;i<100;i++){
7: m += i*x+y;
8: n = z*y;
9: A[i] += A[m];
10:
11: if (debug){
12: printf("m: %d\n",m);
13:}
14: q += i*x;
15:}
16:
17:printf("result:%d %d %d %d %d %d\n",m,n,q,x,y,z);
```

ANSWERS:

constant propagation: line4 changed to z=5+2  $\,$ 

**13.** The following C program is run (with no optimizations) on a processor with a direct-mapped cache that has eight-word (32-byte) blocks and holds 256 words of data: [5m]

```
int i,j,c,stride,array[512];
...
for(i=0; i<10000; i++)
    for (j=0; j<512; j+=stride)
        c += i%(array[j]);</pre>
```

If we consider only the cache activity generated by references to the array and we assume that integers are words, what possible **miss rates** (there may be multiple) can we expect

- 1. if stride = 256?
- 2. if stride = 255?

Explain your answers clearly in a few sentences.

## [OR]