

- ```

Lock: lw R1, mylock # R1 = &mylock
 bne R1, R0, Lock # if (R1 != 0) jump to Lock
 TAS R1, mylock # atomically_do {R1 = &mylock; mylock = 1;}
 bne R1, R0, Lock # if (R1 != 0) jump to Lock
 ret

Unlock: sw 0, mylock # write 0 into &mylock
 ret

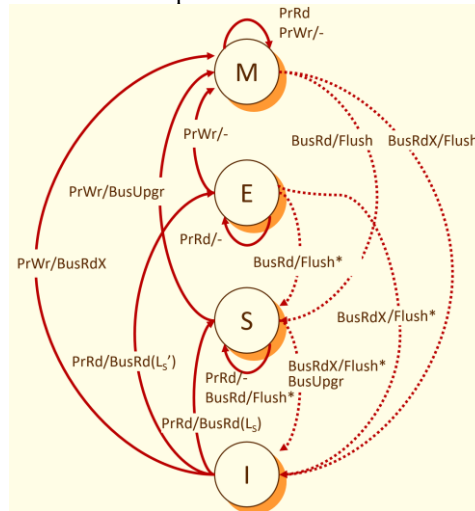
```

By using the following tables, show the operations and bus transactions (or comments): A) in the best case (least number of transactions) and B) in the worst case (highest number of transactions)

[illegible][illegible]

- 2) (POINTS 15/40) Write a CUDA program that reads a color array (`int color[1024]`) and writes an array “`int histogram[256]`” that contains the frequency of each of 256 possible colors (the 256 values are the values that each element of `color[]` can assume). The program should be written in a way that it exploits Thread Level Parallelism as offered by CUDA (a serial or serialized version **has to be avoided**). Hint: try to perform operations in a hierarchical way and use CUDA shared memory.

1) Remembering the state diagram for the MESI protocol:



1A) The best case happens if the interleaving of the operations is such that each processor attempts and get access to the critical section one after the other.

| Bus Trans. Number | Processor Operation | P1 | P2 | P3 | P4 | Bus Transactions/Comments                      |
|-------------------|---------------------|----|----|----|----|------------------------------------------------|
| ---               | (Init.state)        | S  | S  | S  | S  | Initially, P1 holds the lock                   |
| 1                 | sw1                 | M  | I  | I  | I  | <b>BusUpgr</b> – P1 releases the lock          |
| 2                 | lw2                 | S  | S  | I  | I  | <b>BusRd(Ls)/Flush</b> – P2 reads the lock     |
| 3                 | TAS2                | I  | M  | I  | I  | <b>BusUpgr</b> – P2 tries to lock and succeeds |
| --                | sw2                 | I  | M  | I  | I  | P2 releases the lock                           |
| 4                 | lw3                 | I  | S  | S  | I  | <b>BusRd(Ls)/Flush</b> – P3 reads the lock     |
| 5                 | TAS3                | I  | I  | M  | I  | <b>BusUpgr</b> – P3 tries to lock and succeeds |
| --                | sw3                 | I  | I  | M  | I  | P3 releases the lock                           |
| 6                 | lw4                 | I  | I  | S  | S  | <b>BusRd(Ls)/Flush</b> – P4 reads the lock     |
| 7                 | TAS4                | I  | I  | I  | M  | <b>BusUpgr</b> – P4 tries to lock and succeeds |
| --                | sw4                 | I  | I  | I  | M  | P4 releases the lock                           |

1B) The worst case happens if the interleaving of the operations is such that each processor attempts simultaneously the “lw” to read the status of mylock and then simultaneously try to get the access through the TAS instruction.

| Bus Trans. Number | Processor Operation | P1 | P2 | P3 | P4 | Bus Transactions/Comments                   |
|-------------------|---------------------|----|----|----|----|---------------------------------------------|
| ---               | (Init.state)        | S  | S  | S  | S  | Initially, P1 holds the lock                |
| 1                 | sw1                 | M  | I  | I  | I  | <b>BusUpgr</b> -- P1 releases the lock      |
| 2                 | lw2                 | S  | S  | I  | I  | <b>BusRd(Ls)/Flush</b> – P2 reads the lock  |
| 3                 | lw3                 | S  | S  | S  | I  | <b>BusRd(Ls)/Flush*</b> – P3 reads the lock |
| 4                 | lw4                 | S  | S  | S  | S  | <b>BusRd(Ls)/Flush*</b> – P4 reads the lock |
| 5                 | TAS2                | I  | M  | I  | I  | <b>BusUpgr</b> – P2 gets the lock           |
| 6                 | TAS3                | I  | I  | M  | I  | <b>BusRdX/Flush</b> - no lock               |
| 7                 | TAS4                | I  | I  | I  | M  | <b>BusRdX/Flush</b> - no lock               |
| 8                 | st2                 | I  | M  | I  | I  | <b>BusRdX/Flush</b> -- P2 releases the lock |
| 9                 | lw3                 | I  | S  | S  | I  | <b>BusRd(Ls)/Flush</b> – P3 reads the lock  |
| 10                | lw4                 | I  | S  | S  | S  | <b>BusRd(Ls)/Flush*</b> – P4 reads the lock |
| 11                | TAS3                | I  | I  | M  | I  | <b>BusUpgr</b> – P3 gets the lock           |
| 12                | TAS4                | I  | I  | I  | M  | <b>BusRdX/Flush</b> - no lock               |
| 13                | sw3                 | I  | I  | M  | I  | <b>BusRdX/Flush</b> -- P3 releases the lock |
| 14                | lw4                 | I  | I  | S  | S  | <b>BusRd(Ls)/Flush</b> – P4 reads the lock  |
| 15                | TAS4                | I  | I  | I  | M  | <b>BusUpgr</b> – P4 gets the lock           |
| ---               | sw4                 | I  | I  | I  | M  | P4 releases the lock                        |
|                   |                     |    |    |    |    |                                             |
|                   |                     |    |    |    |    |                                             |
|                   |                     |    |    |    |    |                                             |

Note: Depending on the implementation a BusRdX could be directly associated to the (atomic) TAS instruction.

- 2) This is the CUDA code for a possible implementation of the requested kernel (tested on Tesla C1060 with Compute Capability 1.3 and CUDA 4.1):

```
#include <stdio.h>
#include <cuda_runtime.h>
#include <cutil_inline.h>
#include <device_functions.h>
#include <sm_11_atomic_functions.h>

typedef unsigned char uchar;
typedef unsigned int uint;
#define HISTOGRAM_BIN_COUNT 256
#define N 1024
__global__ void histogram3(uint* histogram, uchar* color, int size)
{
 __shared__ uint data[HISTOGRAM_BIN_COUNT];

 // I n i t i a l i z a t i o n
 int stride = blockDim.x;
 for (uint i = threadIdx.x; i < HISTOGRAM_BIN_COUNT; i += stride)
 data[i] = 0;
 __syncthreads();

 // C a l c u l a t e p r i v a t e h i s t o g r a m
 stride = blockDim.x * gridDim.x;
 for (uint i = threadIdx.x + blockDim.x * blockIdx.x;
 i < size; i += stride)
 atomicAdd(&data[color[i]], 1);
 __syncthreads();

 // U p d a t e g l o b a l h i s t o g r a m
 stride = blockDim.x;
 for (uint i = threadIdx.x; i < HISTOGRAM_BIN_COUNT; i += stride)
 atomicAdd(&(histogram[i]), data[i]);
}

int main() {
 uchar* hColor = (uchar*)malloc(N * sizeof(uchar));
 uint* hHistogram3 = (uint*)malloc(HISTOGRAM_BIN_COUNT * sizeof(uint));
 dim3 block, grid;
 uchar* dColor;
 uint* dHistogram;
 cudaMalloc(&dHistogram, HISTOGRAM_BIN_COUNT * sizeof(uint));
 cudaMalloc(&dColor, N * sizeof(uchar));
 srand(2017);
 for (uint i = 0; i < N; ++i) hColor[i] = (uchar)(rand() % 256);
 cudaMemcpy(dColor, hColor, N * sizeof(uchar), cudaMemcpyHostToDevice);
 cudaMemcpy(dHistogram, 0, HISTOGRAM_BIN_COUNT * sizeof(uint));
 block.x = 512;
 grid.x = (N + block.x - 1) / block.x;
 histogram3<<<grid,block>>>(dHistogram, dColor, N);
 cudaMemcpy(hHistogram3, dHistogram,
 HISTOGRAM_BIN_COUNT * sizeof(uint), cudaMemcpyDeviceToHost);
 for (int i = 0; i < HISTOGRAM_BIN_COUNT; ++i)
 printf("%d ", hHistogram3[i]); printf("\n");
}
```

Makefile:

```
EXECUTABLE := histo256
CUFILES_sm_13 := histo256.cu
GENCODE_ARCH := -gencode=arch=compute_13,code=\"sm_13,compute_13\"
 -gencode=arch=compute_20,code=\"sm_20,compute_20\"
include ../../common/common.mk
```