Register number _____

**SRM Institute of Science and Technology**
**College of Engineering and Technology**
**School of Computing**
SRM Nagar, Kattankulathur – 603203, Chengalpattu District, Tamilnadu
**Academic Year: 2023-24 (EVEN)**
**B.Tech-Computer Science & Engineering**          **SET - D**

| | |
|---|---|
| **Test: CLA-T2** | **Date: 28.03.2024** |
| **Course Code & Title: 18CSE419T & GPU Programming** | **Duration: 2 period**s |
| **Year & Sem: III Year /VI Sem** | **Max. Marks: 50** |

**Course articulation matrix:**

| | PO 1 | PO 2 | PO 3 | PO 4 | PO 5 | PO 6 | PO 7 | PO 8 | PO 9 | PO 10 | PO 11 | PO 12 | PSO 1 | PSO 2 | PSO 3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CO-1 | 3 | | | | | | | | | | | | | | 3 |
| CO-2 | | 3 | 2 | | | | | | | | | | | | 3 |
| CO-3 | | 3 | 3 | | | | | | | | | | | | 3 |
| CO-4 | | 3 | 3 | | | | | | | | | | | | 3 |
| CO-5 | | | 3 | 1 | | | | | | | | | | 2 | 3 |

| Part – A(1*10=10 Marks)<br>Answer All the Questions | | | | | |
|---|---|---|---|---|---|
| **Q. N o** | **Questions** | **Mark s** | **B L** | **CO** | **P O** | **PI Cod e** |
| 1 | Calling a kernel is typically referred to as -----------<br>a) Kernel Thread<br>b) Kernel initialization<br>c) Kernel termination<br>d) Kernel invocation | 1 | 2 | CO 3 | 2 & 3 | **4.2.1** |
| 2 | a) A code consisting of GRID ,which runs on GPU consisting of a set of<br>a) 32 thread<br>b) 32 block<br>c) Unit block<br>d) Thread block | 1 | 2 | CO 3 | 2 & 3 | **4.2.1** |
| 3 | The maximum number of threads that can be launched in a specific block is<br>a) 8<br>b) 32<br>c) 256<br>d) 1024 | 1 | 2 | CO 3 | 2 & 3 | **4.2.1** |
| 4 | NVDIA CUDA warp is made up of how many threads?<br>a) 512<br>b) 1024<br>c) 312<br>d) 32 | 1 | 2 | CO 3 | 2 & 3 | **4.2.1** |
| 5 | For a vector addition, assume that the vector length is 2000, each thread calculates one output element and the thread block size is 512 threads. How many threads will be in the grid?<br>a) 2000<br>b) 2024<br>c) 2048<br>d) 2096 | 1 | 2 | CO 3 | 2 & 3 | **4.2.1** |

| 6 | The smallest CUDA thread block dimension is<br>   a) 8<br>   b) 16<br>   c) 32<br>   d) 64 | 1 | 2 | CO3 | 2 & 3 | **4.2.1** |
|---|---|---|---|---|---|---|
| 7 | Which of the following memory locations is common for all the SMs in a typical CUDA GPU?<br>   a) Thread-local memory<br>   b) L1 cache<br>   c) L2 cache<br>   d) Shared memory | 1 | 3 | CO3 | 2 & 3 | **4.2.1** |
| 8 | What is the term used for the combination of CPU and GPU in a hybrid computing system?<br>   a) Homogenous computing<br>   b) Many-core architecture<br>   c) Hardware accelerator<br>   d) Heterogenous computing | 1 | 3 | CO3 | 2 & 3 | **4.2.1** |
| 9 | The scope of a constant memory is<br>   a) Thread<br>   b) Block<br>   c) Warp<br>   d) Grid | 1 | 3 | CO3 | 2 & 3 | **4.2.1** |
| 10 | If each CUDA block can hold a maximum of 512 threads then how many CUDA blocks would be created to process 4000 vector elements<br>   a) 7<br>   b) 8<br>   c) 10<br>   d) 16 | 1 | 3 | CO3 | 2 & 3 | **4.2.1** |

**Part – B (4*4=16 marks)**
**Answer any four Questions**

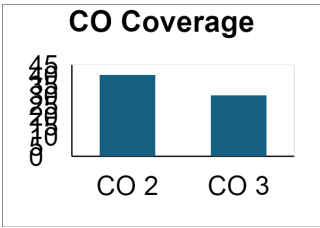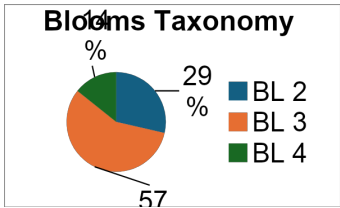| Q.No | Question | Marks | BL | CO | PO | PI Code |
|---|---|---|---|---|---|---|
| 11 | Write a CUDA code to add two numbers and explain the CUDA API functions involved in it. | 4 | 2 | CO2 | 2 | **4.2.1** |
| 12 | What is Texture memory? State the uses of it. | 4 | 3 | CO3 | 2 | **4.2.1** |
| 13 | Name two tools provided by NVIDIA to debug CUDA applications and describe its uses. | 4 | 2 | CO3 | 2 | **4.2.1** |
| 14 | What is zero-copy memory in GPU? What is the need for it? | 4 | 3 | CO3 | 3 | **4.2.1** |
| 15 | List the guidelines for optimizing the performance of the kernel through blocks and grid design. | 4 | 3 | CO3 | 3 | **4.2.1** |

**Part – C (2*12=24 marks)**
**Answer any Two Questions**

| 16 | Brief on the CUDA grid organization as a 3D array of blocks and state the functions involved to access a particular thread. | 12 | 2 | CO2 | 2 | **4.2.1** |
|---|---|---|---|---|---|---|
| 17 | Interpret a static shared memory allocation kernel through a CUDA Histogram calculation. | 12 | 3 | CO3 | 3 | **4.2.1** |
| 18 | Describe the role of constant memory in GPU and recognize how the values are cached and broadcasted using constant memory with an | 12 | 3 | | 3 | **4.2.1** |

| | | | CO 3 | | |
|---|---|---|---|---|---|
| example. | | | | | |

### Blooms Taxonomy



14%
29%
57%

- BL 2
- BL 3
- BL 4

### CO Coverage



45
40
35
30
25
20
15
10
0

CO 2   CO 3

**Approved by Audit Professor/ Course Coordinator**