,	۱		
ŀ	3	L	

# **Question Paper Code: U8G63**

### B.E. / B.Tech. DEGREE EXAMINATION, APRIL / MAY 2025

#### One Credit Course

## Artificial Intelligence & Machine learning

#### 21UAM863 – GPU PROGRAMMING

(Regulations 2021)

Duration: 1.30 hours				Maximum: 50 Marks			
		Answ	er ALL Qu	uestions			
		PART A	- (10 x 1 =	10 Marks	)		
1.	Which GPU architecture introduced CUDA programming?						CO1- U
	a) Kepler	b) Tesla	c) N	Maxwell	ď	) Turing	
2.	2. Threads in CUDA execute in parallel with in						CO1- U
	a) Block	b) Grid	c) (	CPU		d) Kerne	el
3.	Which function is us	sed to specify the	number of	f blocks in	a grid?		CO1- U
	a) blockDim.x	b) gridDim	.x c	) threadIdx	X.X	d) warpSize	
4.	What is the primary challenge of Multi-GPU computing?						CO1- U
	a) Lack of GPU memory b) Synchronization and data transfer overh				transfer overhea	ad	
	c) Slow CPU proces	sing	d) Inability	y to execut	e paralle	el tasks	
5.	Which technology a systems?	allows multiple	GPUs to c	ommunica	ite in N	VIDIA	CO1- U
	a) Hyper-Threading	b) NVLink	c) D	irectX		d) OpenMP	
6.	Problem decomporation programming?	sition is most	useful	in which	type	of	CO1- U
	a) Use shared memo	ory as a cache	b) Use	more CPU	process	ing	
	c) Increase the numb	per of threads	d) Use	more warp	diverge	ence	

7.	The	effect of poor load balancing in parallel programs is:		CO1- U
	a) I	ncreased CPU usage		
	b) I	Decreased execution efficiency due to thread underutilization		
	c) E	nhanced memory access speed		
	d) E	Better management of shared resources		
8.	An	algorithm's performance is most affected by:		CO1- U
	a) T	the number of threads used		
	b) T	he algorithm's time and space complexity		
	c) T	the number of available memory registers		
	d) T	he processor's brand		
9.		ich bus interface is commonly used to connect multiple GPUs system?		CO1- U
	a) S	ATA b) PCIe (Peripheral Component Inte	erconnect Exp	press)
	c) U	JSB d) HDMI		
10.	То	eliminate data dependencies, one approach is	CC	)1- U
	a) U	Ise only sequential execution		
	b) I	ncrease the number of threads		
	c) A	apply synchronization mechanisms like locks and barriers		
	d) R	Reduce the number of CPU cores		
		$PART - B (2 \times 20 = 40 \text{ Marks})$		
11.	(a)	<ul> <li>(i) Given a problem that requires heavy computation (e.g., image processing, deep learning inference), determine whether a GPU or CPU is more suitable and justify your choice.</li> <li>(ii) Implement a basic parallel reduction algorithm (sum reduction) and optimize it using warp-level primitives.</li> </ul>	CO2-App	(10+10)
	(b)	Implement CUDA-based parallelism with graphical processing unit, for solving the parallel matrix multiplication	CO2-App	(20)

12. (a) Describe the concept of inter-process contention in multi-core CO2-App processors. How does contention for shared resources like cache memory or bus bandwidth impact application performance, and how can software optimizations alleviate these issues?

Or

(b) Identify the fault tolerance techniques help avoid errors in CO2-App distributed computing systems? Discuss mechanisms such as redundancy, checkpoints, and replica consistency, and explain how they minimize the impact of system failures.