

## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

Sr. No.	Page No. / Section No. / Clause No.	Tender Description	Query / Clarification / Suggestions from the Bidder	Responses of pre-bid Queries
1	2	Proposed GPUs solution will not be End of Life (EOL) for 5 years from the date of installation	Request to change:- Proposed GPUs solution will not be End of Support Life (EOSL) for 5 years from the date of installation as EOL refer to End of Sale. EOSL ensure proposed solution is covered for comprehensive support from date of installation till it is declared EOSL.	Proposed GPU solution should not be End of Support Life (EOSL) for 5 Years from the date of installation.
2	10	16 Gbps Host Bus Adaptor for connecting with storage.	Request to change:- 16 Gbps Host Bus Adaptor or 25/100G Ethernet adapter for connecting with storage. Alternately, HDR can used for connecting to storage. Justification:- AI works better on Ethernet/Infiniband, hence we don't qualify large GPU nodes with FC connectivity.	Please see corrigendum.
3	10	The solution should be delivered with 1PB usable post RAID 6 or better configuration, expandable up to 2 PB. The proposed storage array should be configured in no single point of failure including controller (at least 2), cache, power supply, cooling fans, etc. Disks: 1) NVMe SSDs and 2) NL-SAS disks. The storage should be distributed with namespace consistent across nodes. IOPS: Min. 10,00,000 Throughput: Min. 100 Gbps Bidirectional	As the total Nodes for AI training: Total Qty-4 sets. Please confirm total Storage requirement is 1PB usable or 4 set of 1PB usable?	1PB Total usable external storage as per specifications
4		The solution should be delivered with 1PB usable post RAID 6 or better configuration, expandable up to 2 PB. The proposed storage array should be configured in no single point of failure including controller (at least 2), cache, power supply, cooling fans, etc. Disks: 1) NVMe SSDs and 2) NL-SAS disks. The storage should be distributed with namespace consistent across nodes. IOPS: Min. 10,00,000 Throughput: Min. 100 Gbps Bidirectional	Request to clarify if this is Block, File or Object storage. Request to consider NVMe SSDs and remove NLSAS as performance requirement is higher. What is the R/W Ratio and Block size for 10,00,000 IOPS? Is throughput 100Gbps or 100GB/s (Gigabytes per second) What will be the frontend interface for storage connectivity?	Please see corrigendum.
5	11	and licenses to be supplied along with product. Full-stack reference designs with all of the leading Storage providers.	Kindly relax the requirement for 'Full-stack reference design', since HPE/OEM will ensure that the proposed GPU nodes and Storage will be compatible in the solution and work seamlessly.	Please see corrigendum.
6	9	Master Node-Generic Query	Bidder need to consider new Infiniband switch or existing SDC switch can be uses if Yes please provide existing switch details	Bidder has to provide new separate said connectivity for completion of propose solution.
7	9	Master Node-Generic Query	No. of master node required ?	7 number Master nodes required for complete propose solution.
8	10	AI Node-Generic Query	Job Scheduler need to consider or not ?	Bidder to provide all items needed to run the AI cluster successfully
9		Networks connectivity	Can GIL advise whether the necessary Infiniband and Ethernet switches are already available at site since we do not see any reference of Switches mention as part of the requirement ?	Bidder has to provide new separate network connectivity for completion of propose solution.
10	11	The Cluster Manager must allow for the easy deployment and management of servers across multiple data centers, the public cloud, and edge locations as a single shared infrastructure through a single interface	Kindly delete this requirement since this is vendor specific and change it to 'The Cluster Manager must allow for the easy deployment and management of servers'	As per RFP.
11	10. Technical Specification: Master Node / 9/15	<b>Network</b> : One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6	<b>Modification</b> : One single Port 200Gb/s HDR InfiniBand Adapter ConnectX-6 <b>Justification</b> : Request you to change to single port HDR as that is most avaiabel with leading OEMs.	Please see corrigendum.
12	10. Technical Specification: Nodes for AI training / 10/15	HBA Card : 16 Gbps Host Bus Adaptor for connecting with storage.	<b>Clarification</b> : Kindly confirm what type of storage connectiuvity is reqd 16Gbps FC or Infiniband.	Please see corrigendum.
13	10. Technical Specification: Nodes for AI training / 10/15	External Storage : Disks: 1) NVMe SSDs and 2) NL-SAS disks.	<b>Modification</b> : Disks: 1) NVMe SSDs / SAS SSD/SATA SSDs and 2) NL-SAS disks. OR The storage must be End to end NVMe. <b>Justification</b> : We request for modification to include optimal technologies to deliver the required performance .	Please see corrigendum.
14	10. Technical Specification: Nodes for AI training / 10/15	IOPS: Min. 10,00,000 Throughput: Min. 100 Gbps Bidirectional	Modification : IOPS Min 6,50,000. Clarification : Please confirm the asked throughput is Giga bits and NOT Giga Bytes!!	Please see corrigendum.

Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)				
Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar				
15	10. Technical Specification: Nodes for AI training / 10/16	<b>System Network</b> : Infiniband for storage delivery	Clarification : Ples remove this clause as Leading storage OEM do not have Infiniband as front end connectivity.	Please see corrigendum.
17	11.Inference Node, Page-11/15	The bidder should configured Head/ Master/ Management Node in (1+1) HA mode deliver the solution.	<b>Clarification</b> : Kindly confirm if the Training and Inference nodes will be in seperate clusters or a single cluster with Master/Head Nodes in High Availability. Also for HA it is highly recommended to have a seperate dedicated shared storage for storing cluster database across the 2 head/Master nodes. Hence we request that a discrete storage array with capacity of about 20TB in raid5 on SSD be asked.	Training and Inference nodes should be in seperate clusters.
18	11.Inference Node, Page-11/16	2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator	<b>Modification</b> : 2 x Accelerators per node, each with minimum. 94GB GPU per Accelerator <b>Justification</b> : Since the ask is for Multi Instance GPU, which is not supported on 40GB GPUs, hence request GIL to amend the clause as required for the solution.	Please see corrigendum.
19	ePBG Percentage(%)	5% of the contract value	As per the recent Amendment in GFR 2017, Performance Security should be 3 to 5% of the contract value. Therefore, we request you to revise the ePBG percentage to 3%.	As per RFP.
20	10. Technical Specification: Master Node / 9/15	<b>Network</b> : One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6	<b>Modification</b> : One single Port 200Gb/s HDR InfiniBand Adapter ConnectX-6 <b>Justification</b> : Request you to change to single port HDR as that is most available with leading OEMs.	Please see corrigendum.
21	10. Technical Specification: Nodes for AI training / 10/15	HBA Card : 16 Gbps Host Bus Adaptor for connecting with storage.	<b>Clarification</b> : Kindly confirm what type of storage connectivity is required 16Gbps FC or Infiniband.	Please see corrigendum.
22	10. Technical Specification: Nodes for AI training / 10/15	External Storage : Disks: 1) NVMe SSDs and 2) NL-SAS disks.	<b>Modification</b> : Disks: 1) NVMe SSDs / SAS SSD/SATA SSDs and 2) NL-SAS disks. OR The storage must be End to end NVMe. <b>Justification</b> : We request for modification to include optimal technologies to deliver the required performance.	Please see corrigendum.
23	10. Technical Specification: Nodes for AI training / 10/15	IOPS: Min. 10,00,000 Throughput: Min. 100 Gbps Bidirectional	Modification : IOPS Min 6,50,000. Clarification : Please confirm the asked throughput is Giga bits and NOT Giga Bytes!!	Please see corrigendum.
24	10. Technical Specification: Nodes for AI training / 10/16	<b>System Network</b> : Infiniband for storage delivery	Clarification : Please remove this clause as Leading storage OEM do not have Infiniband as front end connectivity.	Please see corrigendum.
25	11.Inference Node, Page-11/15 Processors & performance :	The bidder should configured Head/ Master/ Management Node in (1+1) HA mode deliver the solution.	<b>Clarification</b> : Kindly confirm if the Training and Inference nodes will be in seperate clusters or a single cluster with Master/Head Nodes in High Availability. Also for HA it is highly recommended to have a seperate dedicated shared storage for storing cluster database across the 2 head/Master nodes. Hence we request that a discrete storage array with capacity of about 20TB in raid5 on SSD be asked.	Training and Inference nodes should be in seperate clusters.
26	11.Inference Node, Page-11/15 Processors & performance :	2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator	<b>Modification</b> : 2 x Accelerators per node, each with minimum. 94GB GPU per Accelerator <b>Justification</b> : Since the ask is for Multi Instance GPU, which is not supported on 40GB GPUs, hence request GIL to amend the clause as required for the solution.	Please see corrigendum.
27	8. IMPLEMENTATION TIMELINES & PENALTIES	Submission of PBG - Within 15 Days from date of issuance of GEM contract	Based on our experience with various OEM's, We request you to revise the implementation time line as mentioned below Submission of PBG - Within <b>30 Days</b> from date of issuance of GEM contract	As per RFP.
28	8. IMPLEMENTATION TIMELINES & PENALTIES	Supply of the Hardware including Licenses and OEM Warranty Certificate. T1=T+60 days from the date of issuance of contract over GEM	Supply of the Hardware including Licenses and OEM Warranty Certificate. <b>T1=T+90 days</b> from the date of issuance of contract over GEM	Please see corrigendum.
29	8. IMPLEMENTATION TIMELINES & PENALTIES	Installation, commissioning & integration of GPU servers at GSDC along with HLD, LLD documents T2=T1+30	Installation, commissioning & integration of GPU servers at GSDC along with HLD, LLD documents <b>T2=T1+60 days</b>	As per RFP.
30	8. IMPLEMENTATION TIMELINES & PENALTIES	Deployment of required Skilled Resource at GSD T2+7 Days	Deployment of required Skilled Resource at GSD <b>T2+30 Days</b>	As per RFP.
31	b. SLA for Uptime (99.741%)	Uptime of solution <=99.741% In case of failure of proposed solution and non-maintaining targeted value, 0.5% of yearly CAMC payment for every hourly delay in resolution; with max cap of 10 % of total 5 years CAMC value; it will be adjusted from Yearly CAMC payment.	We request you to revise this penalty clause as Uptime of solution <=99.741% In case of failure of proposed solution and non-maintaining targeted value, 0.05% of yearly CAMC payment for every hourly delay in resolution; with max cap of 10 % of Yearly CAMC value; it will be adjusted from Yearly CAMC payment.	As per RFP.
32	Page no. 10, Internal Storage	Item - Nodes for AI Training	1. Kindly Clarify whether 1PB is required on per node basis. Current specifications do not provide clarity about the same	1PB Total usable external storage as per specifications
33	Page no. 10, Internal Storage	Item - Nodes for AI Training	Distribution of storage space among NVMe SSD and NL SAS disks is not specified. Kindly clarify the same	Please consider NVMe.

Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)				
Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar				
34	Page no. 10 , Internal Storage . Item - Nodes for AI Training	Item - Nodes for AI Training	Storage Bandwidth has not been specified , kindly clarify.	Please see corrigendum.
35	Page no. 10 ,	HBA Card .	Kindly confirm whether connectivity between Nodes and Storage is required to be FC based .	Please see corrigendum.
36	Page no. 10 ,	Item - Nodes for AI Training	If FC HBA is required for storage connectivity , then "Infiniband for storage delivery" mentioned is not clear. Kindly confirm.	Please see corrigendum.
37	Page no. 10 ,	Scalability, Cluster and Management Hardware and . Item - Nodes for AI Training	"Full-stack referencedesigns with all of the leading Storage providers." Since full stack reference design is required to be taken from GPU OEM , hence kindly clarify how will Storage Provider provide a full stack design reference . Kindly confirm. If an example can be provided then that will be helpful	Please see corrigendum.
38	Page no. 10	Item - Nodes for AI Training	1. Infiniband for compute communication 2. Infiniband for storage delivery 3. Ethernet (10GbE) for cluster orchestration 4. Ethernet (10GbE) for perimeter connectivity 5. Ethernet (1GbE) for in-band management" In reference to this clause we request clarity about i) size of the interconnect network required for each type of network , ii) blocking ratio of each network , iii) topology required for the networks and iv) scalability provisions required for each network	Please see corrigendum.
39	Page no. 10	Item - Nodes for AI Training	Since the system performance varies with GPU system architecture . Hence in this case if varied architecture based systems are offered by different OEMs in that case unfair comparison shall take place. Hence in order to provide a fair comparison kindly specify the system architecture and the desired performance with <b>benchmarks required to be achieved.</b>	Please see corrigendum.
40	Page no. 10 , System Network	Item - Nodes for AI Training	Since the reference architecture is required to be met as per bid specifications hence it is necessary to provide the details of <b>GPU-GPU - interconnect port ratio without which it will not be possible to offer the required solution</b>	As per RFP.
41	Page no. 10	Item - Nodes for AI Training	Kindly clarify the usage scenario of GPUs in the AI cluster. How the cluster provide GPUs to the users is not defined in the RFP. Hence clarity about required clustering Software is absent. Request you to please provide the same.	Bidders need to ensure that the quoted GPU systems can be provided to end users as - Bare Metal Nodes, Namespaced GPU or GPU partitions
42	Page no. 12 , System Network .	Item - Inference Nodes	1. Ethernet (10GbE) for cluster orchestration 2. Ethernet (10GbE) for perimeter connectivity 3. Ethernet (1GbE) for in-band management" . Since the storage connectivity is absent here , hence request clarity about how will these nodes store and process data to and from.	Please see corrigendum.
43	Page no. 12 , Hypervisor	Item - Inference Nodes	"Hypervisor with Enterprise level highest license and support available should be provided from day one." Since in this clause functionality has not been defined instead a subjective terms called "highest" is specified which varies as a parameter in case of different OEMs hence we request you to please provide clarity about the same	Highest means enterprise support available for 24*7, 365 days with fastest response (Critical support).
44	Page no. 10 , Number of GPUs and GPU Communication.	Item - Nodes for AI Training	"8 x Accelerators per node, each with minimum. 80GB GPU per Accelerator.Minimum 800GB/s bidirectional communication bandwidth perGPU.Should support Tensor core/Matrix core, CUDA / Stream. Processors/ openCL /ROCm with Accelerators," Please note that this clause permits bidders/OEMs to quote existing GPUs which are End of Life but not End of Support and are already succeeded by newer generation GPUs. IN this situation it is requested to your organization to prefer to provide more clarity in order to avoid such participation so that GPU infrastructure offered to you and be sustainable and scalable as well	Please see corrigendum.
45	ATC (GPU Compute for uses AI / ML at GSD)Inference Node: Number of GPUs and GPU Communication	2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator. Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL /ROCm with Accelerators.	The specified is with 40 GB and Nvlink connectivity, considering the current scenario there is no low end card that offer 40 GB of memory with Nvlink connectivity. So request you to let us know if we can offer 1 cards with said memory ?	Please see corrigendum.
46	Regarding Overall Architecture		We request you to please define the cluster orchestration layer details without which the cluster usage is not clear.	Bidders need to ensure that the quoted GPU systems can be provided to end users as - Bare Metal Nodes, Namespaced GPU or GPU partitions
47	Regarding Storage Delivery Network		Kindly specify the design of Storage Delivery Network.	Bidder to design the storage network in such a way that all nodes should be able to talk to the storage network. Storage should deliver the performance mentioned in the specifications.
48	Regarding , Eligibility Condition , Page 1		"The OEM should have executed similar GPU setup for min 3 clients in last 3 Years in India as on date of bid submission. Out of which One client deployment should be One project having similar works total value of INR 100 Cr. Note: Similar works means SITC OF GPU ACCELERATED with multiple GPU Node." We request you to please clarify here that OEM referred here is "OEM for Nodes for AI training" without which there will be an ambiguity in the criteria.	As per RFP.
49	GENRIC	Networking	Could you please confirm if the networking switch will be provided, or if we need to supply it? Additionally, could you clarify the required quantity, switching specifications, and the number of ports needed?	Bidder to provide all items needed to run the AI cluster successfully
50	GENRIC	RACK	Kindly confirm if the rack will be provided, or should we include it in our offering?	Bidders to quote all components including racks required to setup the cluster successfully.

Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)				
Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar				
51	Section 10: System Memory (Page 9)	512 GB DDR4 or higher SDRAM with ECC Advance.	Request clarification if you could consider latest Memroy Platform <b>DDR5</b> Also that will help you to make it standardized across the workload because you already considering DDR5 for other workloads.	Bidder is free to quote DDR4 or higher
52	Section 10: Networking, (Page 9)	1G-Gigabit or higher Ethernet LAN RJ45 ports	Confirm if alternative configurations with <b>Intel Ethernet 800 Series</b> and <b>100G compatibility</b> are acceptable. InfiBand is very OEM specific and other OEM will not able to bid hence request you to consider other OEMs as well.	Please see corrigendum.
53	Section 10: Networking, (Page 9)	• One Single Port 200Gb/s NDR InfiniBand Adapter ConnectX-7 (1xOSFP).	Confirm if alternative configurations with <b>Intel Ethernet 800 Series</b> and <b>100G compatibility</b> are acceptable. InfiBand is very OEM specific and other OEM will not able to bid hence request you to consider other OEMs as well.	Please see corrigendum.
54	Section 10: Networking, (Page 9)	• One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6	Confirm if alternative configurations with <b>Intel Ethernet 800 Series</b> and <b>100G compatibility</b> are acceptable. InfiBand is very OEM specific and other OEM will not able to bid hence request you to consider other OEMs as well.	Please see corrigendum.
55	Section 10: Networking, (Page 10)	Minimum 4 nos. of InfiniBand NDR ports for internode communication, 1 nos. of port for BMC (dedicated LAN port), minimum 1 no. of 1 GbE & 10 GbE (Fiber) port. Additional InfiniBand HDR/NDR200 single/2-port HCA for storage connectivity.	InfiBand is very OEM specific and other OEM will not able to bid hence request you to consider other OEMs as well.	Please see corrigendum.
56	System Network (Page 10)	Following N/W are required: 1. Infiniband for compute communication 2. Infiniband for storage delivery 3. Ethernet (10GbE) for cluster orchestration 4. Ethernet (10GbE) for perimeter connectivity 5. Ethernet (1GbE) for in-band management	InfiBand is very OEM specific and other OEM will not able to bid hence request you to consider other OEMs as well.	Please see corrigendum.
57	System Network (Page 11 & Page 12)	Following N/W are required: 1. Ethernet (10GbE) for cluster orchestration 2. Ethernet (10GbE) for perimeter connectivity 3. Ethernet (1GbE) for in-band management	InfiBand is very OEM specific and other OEM will not able to bid hence request you to consider other OEMs as well.	Please see corrigendum.
58	OS Support	The system should support latest version of Red Hat Enterprise Linux / Ubuntu Linux server Quoted OS should be under Enterprise support from OEM.	The system should support latest version of Red Hat Enterprise Linux / Ubuntu Linux / <b>RHEL AI / Red Hat OpenShift AI server</b> Quoted OS should be under Enterprise support from OEM with premium or highest level of support.	As per RFP.
59	Hypervisor	Hypervisor with Enterprise level highest license and support available should be provided from day one.	Hypervisor with Enterprise level highest license and support available should be provided from day one and must have a native open source LLM (large Language model) bundle in the proposed hypervisor. Since DIT/GSDC is planning to have inferencing / training on the proposed hardware as per Specs.	As per RFP.
60	AI Enterprise Software	AI Enterprise software & subscription or equivalent for each and every GPUs to be included from day one of Installation. Software Support All necessary and required software, SDK, libraries, tools to cater and run the AI/ML workload should be provide from day 1.	AI Enterprise software & subscription or equivalent for each and every GPUs to be included from day one of Installation. Software Support All necessary and required software, SDK, libraries, tools to cater and run the AI/ML workload should be provide from day 1. <b>bias and drift detection, better access to accelerators, and a centralized registry to share, deploy, and track models. the proposed solution can be open-source with enterprise support.</b>	As per RFP.
61	MLOps	New Specs to be Added.	the Proposed solution must have a MLOps, Pipeline for the project with native opensource workspace like Jupyter notebook, Pytorch Apache Spark etc. Tight integration with opensource supported CI/CD tools must be provided that will allows ML models to be quickly deployed iteratively, as needed by Govt entity / Dept.	As per RFP.
62	EMD: Page No. 8	1) EMD with Account Payee Demand Draft 2) EMD with Payment online through RTGS / internet banking	We would like to confirm that as per the GEM Standard Terms and Conditions of "BID SECURITY" – Pg. No. 4, Reference File Name: General Terms and Conditions on GeM 4.0 (Version 1.20) dated 16th December 2024 (attached herewith), we intend to submit the Earnest Money Deposit (EMD) in the form of a Bank Guarantee. Additionally, GEM, by default, provides a standard format for the Bank Guarantee for every bid.  In light of this, we kindly request your approval to submit the EMD in the form of a Bank Guarantee, which is widely accepted across all Central/State Governments, PSUs, and other government organizations.	As per RFP.
63	Eligibility Conditions: Clause 4, page no. 2	Bidder's experience	As per standard paractice, Experience is mostly asked as "Bidder / OEM should have.....", we request Authority to kindly amend the Experience clause accordingly.	Please see corrigendum.
64	Eligibility Conditions: Clause 4, page no. 2	The bidder should have experience of set up of similar GPU base solution with minimum 75 GPU/CPU based solution with minimum 2000 cores in last 3 Years in last 3 Years in India with below criteria as on last date of bid submission: - One project having similar works total value of INR 32 Cr or - Two project having similar works total value of INR 20 Cr or - Three project having similar works to-tal value of INR 12 Cr  <b>Note:</b> Similar works means SITC OF GPU AC-CELERATED with multiple GPU Node.	Since execution / deployment of GPU Based AI/ML projects in INDIA is relatively new , the vendors having similar experience may have difficulties in the single/multiple PO value that is 32 Cr / 20 Cr / 12 Cr.  For successful execution of similar project, importantly, the bidder should have exeperice on deploying reasonable size of GPU/CPU based set-up, which is technical skills, hence the commercial (value of project) should not matter.  We request you to kindly asked cumulative experience of execution of similar GPU based projects in consolidation on total number of GPU/CPU based deployments that is cumulative 50 GPUs / 1500 CPUs in <b>past 5 years</b> .  <b>We believe that experience of deployment of GPU/CPU based solution is vital as compared to the value of the project.</b>  We hereby request you to kindly emend the clause as below.  The bidder should have <b>cumulative</b> experience of setting up similar GPU base solution with minimum <b>50 GPU/CPU</b> based solution with minimum <b>1500 cores in last 5 Years</b> in India as on last date of bid submission:  Note: Similar works means SITC OF GPU AC-CELERATED with multiple GPU Node.	Please see corrigendum.
65	Implementation Timelines	Supply of the Hardware	The timeline given for delivery is difficult, kindly make it at least 120 days.	Please see corrigendum.

## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

66	Nodes for AI training Clause 2, page No. 11	8 x Accelerators per node, each with minimum. 80GB GPU per Accelerator. Minimum 800GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators,	The current generation of GPUs provide far more memory than the previous generation. Since this is a new procurement it is our suggestion that GSDC should go with the newer GPU generation with better capacity as it would better capabilities and provide investment protection. It is advised to clarify that the system should be able to support any GPU to GPU communication at 800GB/s bi-directionally and 800GB/s should not be cumulative bandwidth across all GPUs  Request to amend the clause as  "8 x Accelerators per node, each with minimum. <b>141GB GPU per Accelerator.</b> <b>Minimum 800GB/s bidirectional communication bandwidth between any GPU-GPU pair.</b> Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators"	Please see corrigendum.
67	Nodes for AI training Clause 5, page 11	Minimum 4 nos. of InfiniBand NDR ports for internode communication, 1 nos. of port for BMC (dedicated LAN port), minimum 1 no. of 1 GbE & 10 GbE (Fiber) port. Additional InfiniBand HDR/NDR200 single/2-port HCA for storage connectivity.	It is recommended that the ratio between GPU to Port should be 1:1 and each GPU should have its own dedicated port for external communication to avoid congestion and bottlenecks at the time of training.  Request to amend the clause as  "Minimum <b>8</b> nos. of InfiniBand <b>400G</b> NDR ports for internode <b>GPU-to-GPU</b> communication, 1 nos. of port for BMC (dedicated LAN port), minimum 1 no. of 1 GbE <b>or</b> 10 GbE (Fiber) port. Additional InfiniBand NDR single 2-port HCA for storage and user connectivity"	Please see corrigendum.
68	Nodes for AI training External Storage Page No.11	The solution should be delivered with 1PB usable post RAID 6 or better configuration, expandable up to 2 PB. The proposed storage array should be configured in no single point of failure including controller (at least 2), cache, power supply, cooling fans, etc. Disks: 1) NVMe SSDs and 2) NL-SAS disks. The storage should be distributed with namespace consistent across nodes. IOPS: Min. 10,00,000 Throughput: Min. 100 Gbps Bidirectional	Kindly clarify that the external storage should be PFS/NFS based as the GPU based system require PFS/NFS based storage systems. We would like to suggest to get NVMe SSD disks and recommended throughput as min 1GB/s per GPU for good performance. The solution <b>should be PFS/NFS based</b> and delivered with 1PB usable post RAID 6 or better configuration, expandable up to 2 PB. The proposed storage array should be configured in no single point of failure including controller (at least 2), cache, power supply, cooling fans, etc. Disks: 1) <b>NVMe SSDs only</b> . The storage should be distributed with namespace consistent across nodes. IOPS: Min. 10,00,000 Throughput: <b>Min 1GB/s per GPU</b>	Please see corrigendum.
69	Nodes for AI training, Page No 12	Following N/W are required: 1. Infiniband for compute communication 2. Infiniband for storage delivery 3. Ethernet (10GbE) for cluster orchestration 4. Ethernet (10GbE) for perimeter connectivity 5. Ethernet (1GbE) for in-band management	We would suggested to go with 25GbE network for mgmt, orchestration and perimeter connectivity Following N/W are required: 1. NDR based Infiniband for compute communication 2. NDR Infiniband for storage delivery 3. Ethernet (25GbE) for cluster orchestration 4. Ethernet (25GbE) for perimeter connectivity 5. Ethernet (1GbE) for in-band management	Please see corrigendum.
70	Nodes for AI training - Number of GPUs and GPU Communication, page no 12	2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator. Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators.	The current generation of GPUs provide far more memory than the previous generation. Since this is a new procurement it is our suggestion that Authority should go with the newer GPU generation with better capacity as it would better capabilities and provide investment protection.  2 x Accelerators per node, each with minimum. 141GB GPU per Accelerator. Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators.	Please see corrigendum.
71	Nodes for AI training - System Network, Page no 12	Following N/W are required: 1. Ethernet (10GbE) for cluster orchestration 2. Ethernet (10GbE) for perimeter connectivity 3. Ethernet (1GbE) for in-band management	We would to go with 25GbE network for inferencing capacity to support high request periods Following N/W are required: 1. Ethernet (25GbE) for cluster orchestration 2. Ethernet (25GbE) for perimeter connectivity 3. Ethernet (1GbE) for in-band management	Please see corrigendum.
72	Nodes for AI training AI Enterprise Software, page 12	AI Enterprise software & subscription or equivalent for each and every GPUs to be included from day one.	Kindly include that OEM SLA based enterprise support should be available for not just the OEM hardware but OEM software as well  AI Enterprise software & subscription or equivalent for each and every GPUs to be included from day one of Installation. Bidder to ensure that enterprise level OEM support & SLA is available for all OEM provided software and libraries	As per RFP.

## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

73		NEW CLAUSE : Request to add	<p>MLperf is an industry standard benchmark for evaluating GPU based AI systems and platforms and is universally accepted.</p> <p>More than 50 leading industry leaders are members:  <a href="https://mlcommons.org/our-members/">https://mlcommons.org/our-members/</a></p> <p>It is critical that GSDC evaluate all quoted GPU systems as per the mlperf standards  Bidder needs to submit proof of the quoted GPU being listed in MLcommons inferencing benchmarks at the time of bidding</p>	Please see corrigendum.
			<p><b>Image segmentation (medical)</b>  3D-Unet-99  Throughput for single NODE inference (99% offline) = 45 Samples/s or higher</p>	Please see corrigendum.
			<p><b>NLP</b>  Bert-99  Throughput for single NODE inference (99% offline) = 70000 Samples/s or higher</p>	Please see corrigendum.
			<p><b>Recommendation</b>  dlrm-v2-99  Throughput for single NODE inference (99% offline) = 610000 Samples/s or higher</p>	Please see corrigendum.
			<p><b>LLM Summarization</b>  gptj-99  Throughput for single NODE inference (99% offline) = 17000 Tokens/s or higher</p>	Please see corrigendum.
			<p><b>LLM Q&amp;A</b>  LLama2-70B  Throughput for single NODE inference (99% offline) = 27000 Tokens/s or higher</p>	Please see corrigendum.
			<p><b>LLM Text Generation (Math &amp; Code)</b>  Mixtral-8x7B  Throughput for single NODE inference (99% offline) = 53000 Tokens/s or higher</p>	Please see corrigendum.
			<p><b>Image Classification</b>  ResNet  Throughput for single NODE inference (99% offline) = 716000 Samples/s or higher</p>	Please see corrigendum.
			<p><b>Object Detection</b>  RetinaNet  Throughput for single NODE inference (99% offline) = 11400 Samples/s or higher</p>	Please see corrigendum.
			<p><b>Image Generation</b>  Stable Diffusion-XL  Throughput for single NODE inference (99% offline) = 14 Samples/s or higher</p>	Please see corrigendum.
74	10. Technical Specification: Master Node, Point 15, page 9	<b>Network</b> : One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6	<p><b>We Request you to change to single port HDR as that is most available with leading OEMs.</b>  <b>Modification</b> : One single Port 200Gb/s HDR InfiniBand Adapter ConnectX-6</p>	Please see corrigendum.
75	10. Technical Specification: Nodes for AI training, point 15, page 10	HBA Card : 16 Gbps Host Bus Adaptor for connecting with storage.	We request you to kindly confirm what type of storage connectivity is reqd 16Gbps FC or Infiniband.	Please see corrigendum.
76	10. Technical Specification: Nodes for AI training, point 15, page 10	External Storage : Disks: 1) NVMe SSDs and 2) NL-SAS disks.	<p><b>We request for modification to include optimal technologies to deliver the required performance .</b>  Request for Change : Disks: 1) NVMe SSDs / SAS SSD/SATA SSDs and 2) NL-SAS disks. OR The storage must be End to end NVMe.</p>	Please see corrigendum.
77	10. Technical Specification: Nodes for AI training, point 15, page 10	IOPS: Min. 10,00,000 Throughput: Min. 100 Gbps Bidirectional	Request for change : IOPS Min 6,50,000. Clarification : Please confirm the asked throughput is Giga bits and NOT Giga Bytes!!	Please see corrigendum.
78	10. Technical Specification: Nodes for AI training , point 16, page 10	<b>System Network</b> : Infiniband for storage delivery	We request you to kindly remove this clause as Leading storage OEM do not have Infiniband as front end connectivity.	Please see corrigendum.
79	11.Inference Node, point 15, page 11	The bidder should configured Head/ Master/ Management Node in (1+1) HA mode deliver the solution.	We request you to kindly confirm if the Training and Inference nodes will be in separate clusters or a single cluster with Master/Head Nodes in High Availability. Also for HA it is highly recommended to have a separate dedicated shared storage for storing cluster database across the 2 head/Master nodes. Hence we request that a discrete storage array with capacity of about 20TB in raid5 on SSD be asked.	Training and Inference nodes should be in separate clusters.
80	11.Inference Node, point 16, page 11	2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator	<p><b>Since the ask is for Multi Instance GPU, which is not supported on 40GB GPUs , hence request GIL to amend the clause as required for the solution.</b>  <b>Modification</b> : 2 x Accelerators per node, each with minimum. 94GB GPU per Accelerator</p>	Please see corrigendum.

## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

81	Page 9 Master Node:	Processors - Latest Generation Intel® Xeon® platinum or AMD Epyc scalable processors with Minimum Dual 56-Core	Suggested modification to the clause as below -  Processors - 5th Gen Intel® Xeon® platinum or 4th Gen AMD Epyc processors with Dual 64-Core 2.2 GHz minimum. <b>Justification:</b> Intel 5th Gen offers 64 core processors hence request you to consider 64 core processors instead of 56 core processors. And Intel 5th Generation is comparable to AMD EPYC 4th Generation.  Also, both Intel and AMD have launched their next generation of processor series, hence request you to please remove the word 'latest generation'.	As per RFP.
82	Page 9 Master Node:	Network - • 1G-Gigabit or higher Ethernet LAN RJ45 ports • One Single Port 200Gb/s NDR InfiniBand Adapter ConnectX-7 (1x OSFP). • One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6	Suggested modification to the clause as below -  Network - • 1G-Gigabit or higher Ethernet LAN RJ45 ports • One Single Port 200Gb/s NDR InfiniBand Adapter ConnectX-7 (1x OSFP) or Single port 200 Gb Ethernet Port • One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6 or Single port 200 Gb Ethernet Port <b>Justification:</b> Request you to please consider standard ethernet for LAN connectivity. Infiniband solution is available only from a single OEM vendor and proprietary whereas ethernet and FC technology is available from multiple vendors.	Please see corrigendum.
83	Page 10 Nodes for AI training: Total Qty-4 sets.	Processors & performance (per node, minimum) - Min Dual 56-core 4th Gen Intel® Xeon® platinum or AMD Epyc scalable processors, with Min 8 X GPU Accelerators	Suggested modification to the clause as below -  Processors & performance (per node, minimum) - Min Dual 64-core 2.2 GHz 5th Gen Intel® Xeon® platinum or 4th Gen AMD Epyc processors, with Min 8 X GPU Accelerators <b>Justification:</b> For GPU servers, higher no. of cores and CPU speed is very important to get maximum performance out of the GPU servers. Intel 5th Gen offers 64 core processors hence request you to consider 64 core processors instead of 56 core processors.	As per RFP.
84	Page 10 Nodes for AI training: Total Qty-4 sets.	System Memory - The system should be configured with Minimum 2TB DDR5 RAM with all slots populated in balanced configuration for maximum bandwidth.	Suggested modification to the clause as below - System Memory - The system should be configured with Minimum 2TB DDR5 RAM with all Memory slots to be populated considering only 1 Memory DIMM per Memory Channel (1DPC) in balanced configuration for maximum memory bandwidth. <b>Justification:</b> Memory bandwidth/throughput is directly proportionate to runtime memory speed and no. of memory channels populated. This is very important consideration from application performance perspective. Any reduction in memory bandwidth would have major impact on overall server performance and due to the same reason, all the SPEC benchmarks by the server OEMs are done only with 1 memory DIMM populated per memory channel. Please Refer page no. 18 and 19 <a href="https://www.hpe.com/psnow/doc/a50010242enw">https://www.hpe.com/psnow/doc/a50010242enw</a> – Server Memory population rules for Intel Xeon Servers. Please find below the server memory DIMM slots and memory bandwidth calculations for your ready reference - 1. Intel has 8 Memory channels and 2 DIMMs per channel (2 DPC), hence the Intel dual processor server has (2 processors x 8 Memory channels x 2 DPC) = 32 Memory DIMM Slots Intel 8480+ : Intel memory bandwidth with 1 DPC is 8*4800*8/1000 = 307.2 GB/s and with 2 DPC 8*4400*8/1000 = 281.6 GB/s per processor Intel 8570 : Intel memory bandwidth with 1 DPC is 8*5600*8/1000 = 358.4 GB/s and with 2 DPC 8*4400*8/1000 = 281.6 GB/s per processor 2. AMD has 12 Memory channels and 1 DIMMs per channel (1 DPC), hence the AMD dual processor server has (2 processors x 12 Memory channels x 1 DPC) = 24 Memory DIMM Slots AMD 9534 : AMD memory bandwidth with 1 DPC is 12*4800*8/1000 = 460.8 GB/s per processor Hence, it is recommended to populate the memory slots with only 1 DPC for achieving highest server performance.	As per RFP.
85	Page 10 Nodes for AI training: Total Qty-4 sets.	Network - Minimum 4 nos. of InfiniBand NDR ports for internode communication, 1 nos. of port for BMC (dedicated LAN port), minimum 1 no. of 1 GbE & 10 GbE (Fiber) port. Additional InfiniBand HDR/NDR200 single/2-port HCA for storage connectivity.	Suggested modification to the clause as below -  Network - Minimum 4 nos. of InfiniBand NDR ports or 4 nos. of 200 Gb Ethernet Ports for internode communication, 1 nos. of port for BMC (dedicated LAN port), minimum 1 no. of 1 GbE & 10 GbE (Fiber) port. Additional InfiniBand HDR/NDR200 single/2-port HCA or dual port 16 Gb FC HBA for storage connectivity. <b>Justification:</b> Request you to please consider standard ethernet and FC HBA for LAN and SAN connectivity. Infiniband solution is available only from a single OEM vendor and proprietary whereas ethernet and FC technology is available from multiple vendors.	Please see corrigendum.

## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

86	Page 10 and 11 Nodes for AI training: Total Qty-4 sets.	Scalability, Cluster and Management Hardware and software - System should be scalable with multi node cluster. Software support & cluster tools and management hardware and software and licenses to be supplied along with product. Full-stack reference designs with all of the leading Storage providers.	Suggested modification to the clause as below -  Scalability, Cluster and Management Hardware and software - System should be scalable with multi node cluster. Software support & cluster tools to be supplied along with product. <b>Justification:</b> Full stack reference design is available only from a single OEM vendor which needs to be followed exactly as it is which make the solution OEM Vendor Lock-in and does not allow any flexibility. It is suggested to keep the specifications in a neutral way to allow multiple vendors to participate and offer the most cost effective and efficient solution.	Please see corrigendum.
87	Page 11 Nodes for AI training: Total Qty-4 sets.	Certification - Rack Servers should be certified by GPU Controller / Accelerator OEM, the Certificate or listing of offered Server model in GPU Controller / Accelerator OEM website must be submitted along with bid.	Suggested modification to the clause as below -  Certification - Rack Servers should be certified for GPU Controller / Accelerator OEM. The proposed server model datasheet should list the offered GPU Controller / Accelerator as supported. <b>Justification:</b> This is another way of vendor lock-in by controlling all the solution aspects. Hence, it is suggested that the GPU solution should be accepted if the offered server datasheet lists the same.	As per RFP.
88	Page 11 Inference Node: Total Qty-12 sets.	Processors & performance (per node, minimum)  Min Dual 32-core 4th Gen Intel® Xeon® platinum or AMD Epyc scalable processors, with Min 2 X GPU Accelerator. The bidder should configured Head/ Master/ Management Node in (1+1) HA mode deliver the solution.	Suggested modification to the clause as below -  Min Dual 32-core 5th Gen Intel® Xeon® platinum or 4th Gen AMD Epyc scalable processors, with Min 2 X GPU Accelerator. The bidder should configured Head/ Master/ Management Node in (1+1) HA mode deliver the solution. <b>Justification:</b> Intel 5th Gen is already available and is comparable with ADM 4th Gen processors. Hence request you to consider the same.	As per RFP.
89	Page 11 Inference Node: Total Qty-12 sets.	Number of GPUs and GPU Communication - 2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator. Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators.	Suggested modification to the clause as below -  2 x Accelerators per node, each with minimum. 60GB GPU per Accelerator. Minimum 300GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators. <b>Justification:</b> Higher GPU Memory capacity of 60 GB or more recommended. Request you to please consider 300 GB/s GPU-GPU bandwidth as with 2 GPUs AMD MI210 supports 300 GB/s bandwidth.  AMD MI210 Datasheet - <a href="https://www.amd.com/content/dam/amd/en/documents/instinct-business-docs/product-briefs/instinct-mi210-brochure.pdf">https://www.amd.com/content/dam/amd/en/documents/instinct-business-docs/product-briefs/instinct-mi210-brochure.pdf</a>	Please see corrigendum.
90	Page 11 Inference Node: Total Qty-12 sets.	Multi Instance GPU -  Capability to support partitioning of single GPU into multiple GPU instances.	Suggested modification to the clause as below -  Multi Instance GPU - Optional Capability to support partitioning of single GPU into multiple GPU instances. <b>Justification:</b> AMD MI210 GPU does not support GPU partitioning, hence request you to please keep this feature as optional	Please see corrigendum.
91	Page 11 Inference Node: Total Qty-12 sets.	System Memory - The system should be configured with Minimum 2TB DDR5 RAM with all slots populated in balanced configuration for maximum bandwidth.	Suggested modification to the clause as below - System Memory - The system should be configured with Minimum 2TB DDR5 RAM with all Memory slots to be populated considering only 1 Memory DIMM per Memory Channel (1DPC) in balanced configuration for maximum memory bandwidth. <b>Justification:</b> Memory bandwidth/throughput is directly proportionate to runtime memory speed and no. of memory channels populated. This is very important consideration from application performance perspective. Any reduction in memory bandwidth would have major impact on overall server performance and due to the same reason, all the SPEC benchmarks by the server OEMs are done only with 1 memory DIMM populated per memory channel. Please Refer page no. 18 and 19 <a href="https://www.hpe.com/psnow/doc/a50010242enw">https://www.hpe.com/psnow/doc/a50010242enw</a> – Server Memory population rules for Intel Xeon Servers. Please find below the server memory DIMM slots and memory bandwidth calculations for your ready reference - 1. Intel has 8 Memory channels and 2 DIMMs per channel (2 DPC), hence the Intel dual processor server has (2 processors x 8 Memory channels x 2 DPC) = 32 Memory DIMM Slots Intel 8462Y+ : Intel memory bandwidth with 1 DPC is 8*4800*8/1000 = 307.2 GB/s and with 2 DPC 8*4400*8/1000 = 281.6 GB/s per processor Intel 8562Y+ : Intel memory bandwidth with 1 DPC is 8*5600*8/1000 = 358.4 GB/s and with 2 DPC 8*4400*8/1000 = 281.6 GB/s per processor 2. AMD has 12 Memory channels and 1 DIMMs per channel (1 DPC), hence the AMD dual processor server has (2 processors x 12 Memory channels x 1 DPC) = 24 Memory DIMM Slots AMD 9534 : AMD memory bandwidth with 1 DPC is 12*4800*8/1000 = 460.8 GB/s per processor. Hence, it is recommended to populate the memory slots with only 1 DPC for achieving highest server performance.	As per RFP.



## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

92	Page 12 Inference Node: Total Qty-12 sets.	Virtual GPU - Support for virtual GPU to share a physical GPU across multiple VMs. required license should be from day one.	Suggested modification to the clause as below -  Virtual GPU - Optional Support for virtual GPU to share a physical GPU across multiple VMs. required license should be from day one. <b>Justification:</b> AMD MI210 GPU does not support GPU Virtual GPU, hence request you to please keep this feature as optional	Please see corrigendum.
93	Page 12 Inference Node: Total Qty-12 sets.	Certification - Rack Servers should be certified by GPU Controller / Accelerator OEM, the Certificate or listing of offered Server model in GPU Controller / Accelerator OEM website must be submitted along with bid.	Suggested modification to the clause as below -  Certification - Rack Servers should be certified for GPU Controller / Accelerator OEM. The proposed server model datasheet should list the offered GPU Controller / Accelerator as supported. <b>Justification:</b> This is another way of vendor lock-in by controlling all the solution aspects. Hence, it is suggested that the GPU solution should be accepted if the offered server datasheet lists the same.	As per RFP.
94	Page 12 Inference Node: Total Qty-12 sets.	Cluster Management & Scheduler and hardware - Cluster management for system provisioning and monitoring needs to be included, The Cluster Manager must support multiple of hardware vendors. The Cluster Manager must allow for the easy deployment and management of servers across multiple data centers, the public cloud, and edge locations as a single shared infrastructure through a single interface. All necessary hardware, software and necessary licenses should be provided from day 1	Suggested modification to the clause as below -  Cluster Management & Scheduler and hardware - Cluster management for system provisioning and monitoring needs to be included, The Cluster Manager must support multiple of hardware vendors. The Cluster Manager must allow for the easy deployment and management of servers through a single interface. All necessary hardware, software and necessary licenses should be provided from day 1. <b>Justification:</b> Request you to please consider Cluster management capabilities with-in the boundries of the Datacenter.	As per RFP.
95	2025116122226169/ page - 11	Min Dual 32-core 4th Gen Intel® Xeon® platinum or AMD Epyc,scalable processors, with Min 2 X GPU Accelerator. The bidder should configured Head/ Master/ Management Node in (1+1) HA mode deliver the solution.	Clause to be amended to - Min Dual 32-core 4th Gen Intel® Xeon® platinum or AMD Epyc scalable processors with Min 2 X GPU Accelerator or <b>equivalent IBM Power Processor.</b> The bidder should configured Head/ Master/ Management Node in (1+1) HA mode deliver the solution.	Please see corrigendum.
96	2025116122226169/ page - 11	2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator. Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL /ROCm with Accelerators.	Clause to be amended to - <b>2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator or Equivalent MMA Matrix Accelerators</b>	Please see corrigendum.
97	Internal Storage	• For Operating System: Two numbers of 3.84 TB capacity hot swap Enterprise NVMe SSDs with MTBF of 2 million hours or more, Configure in RAID 1. • For Data: Six numbers of 3.84 TB capacity interface hot swap Enterprise NVMe SSDs with MTBF of 2 million hours or more, configured in RAID 5.	please revise this clause as below: "• For Operating System: Two numbers of 3.84 TB capacity hot swap Enterprise NVMe SSDs, Configure in RAID 1. • For Data: Six numbers of 3.84 TB capacity interface hot swap Enterprise NVMe SSDs configured in RAID 5.	Please see corrigendum.
98	Storage Controller	Hardware RAID 0,1, 5, 6, 10, 50, 60 with 8GB cache Flash based cache protection module should be included, should support Gen4.0 PCIe NVMe.	please revise this clause as below: "Hardware RAID 0,1, 5, 6, 10, 50, 60 with min 4GB cache Flash based cache protection module should be included, should support Gen4.0 PCIe NVMe."	Please see corrigendum.
99	Network	• 1G-Gigabit or higher Ethernet LAN RJ45 ports • One Single Port 200Gb/s NDR InfiniBand Adapter ConnectX-7 (1x OSFP). • One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6	• 1G/10G/25G-Gigabit or higher Ethernet LAN RJ45 ports / SFP Ports • One Single / dual Port 200Gb/s NDR InfiniBand Adapter ConnectX-7 (1x OSFP)/Ethernet adapter • One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX- 6/Ethernet adapter	Please see corrigendum.
100	Number of GPUs and GPU Communication Nodes for AI training	8 x Accelerators per node, each with minimum. 80GB GPU per Accelerator. Minimum 800GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL /ROCm with Accelerators,	Please revise as below for wider OEM participation: "8 x Accelerators per node, each with minimum. 100GB GPU per Accelerator or higher. Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL /ROCm with Accelerators,"	Please see corrigendum.
101	Internal Storage Nodes for AI training	• For Operating System: minimum 1.92 TB M.2 NVMe drives • For Data: Minimum 4 * 3.84 TB U.2 NVMe drives	Please revise as below for wider OEM participation: "• For Operating System: minimum 1 TB M.2 NVMe drives • For Data: Minimum 16* 1.9 TB NVMe drives"	Please see corrigendum.
102	HBA Card Nodes for AI training	16 Gbps Host Bus Adaptor for connecting with storage.	Please remove this clause as normally DPUs are required for storage traffic communication	Please see corrigendum.

## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

103	System Network Nodes for AI training	Following N/W are required: 1. Infiniband for compute communication 2. Infiniband for storage delivery 3. Ethernet (10GbE) for cluster orchestration 4. Ethernet (10GbE) for perimeter connectivity 5. Ethernet (1GbE) for in-band management	Please revise as below: "Following N/W are required: 1. Infiniband / Ethernet for compute communication 2. Infiniband / Ethernet for storage delivery 3. Ethernet (25GbE) for cluster orchestration 4. Ethernet (25GbE) for perimeter connectivity 5. Ethernet (1GbE) for in-band management"	Please see corrigendum.
104	Cluster Management & Scheduler and hardware Nodes for AI training	Cluster management for system provisioning and monitoring needs to be included, The Cluster Manager must support multiple of hardware vendors. The Cluster Manager must allow for the easy deployment and management of servers across multiple data centers, the public cloud, and edge locations as a single shared infrastructure through a single interface. All necessary hardware, software and necessary licenses should be provided from day 1	Please revise as below: Cluster management for system provisioning and monitoring needs to be included, The Cluster Manager must support multiple of hardware vendors. The Cluster Manager must allow for the easy deployment and management of servers across multiple data centers, the public cloud, and edge locations as a single shared infrastructure through a single interface. All necessary hardware, software and necessary licenses should be provided from day 1. Cluster Management Software should be Cloud-Native. It should have Multi-Cluster Support From Day one. It should be able to support Multiple GPU OEM. It should have Observability For Infrastructure, Applications available Realtime providing Per GPU, Per Node and Entire Cluster Real-Time Utilization and Suggest required Optimization. Cluster Manager should support low latency Network Deployment Architecture. Cluster Manager Deployment should be able to support Network, Orchestration, Service Mesh, Observability and Security from Single UI Interface.	Please see corrigendum.
105	Inference Node Number of GPUs and GPU Communication	2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator. Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators.	Please revise as below for wider OEM participation: "8 x Accelerators per node, each with minimum. 100GB GPU per Accelerator or higher . Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators,"	Please see corrigendum.
106	10. Technical Specification: Master Node / 9/15	<b>Network</b> : One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6	<b>Modification</b> : One single Port 200Gb/s HDR InfiniBand Adapter ConnectX-6 <b>Justification</b> : Request you to change to single port HDR as that is most available with leading OEMs.	Please see corrigendum.
107	10. Technical Specification: Nodes for AI training / 10/15	HBA Card : 16 Gbps Host Bus Adaptor for connecting with storage.	<b>Clarification</b> : Kindly confirm what type of storage connectivity is reqd 16Gbps FC or InfiniBand. OR we recommend to look at 100Gbps front end GbE	Please see corrigendum.
108	10. Technical Specification: Nodes for AI training / 10/15	External Storage : Disks: 1) NVMe SSDs and 2) NL-SAS disks.	<b>Modification</b> : Disks: 1) NVMe SSDs / SAS SSD/SATA SSDs and 2) NL-SAS disks. OR The storage must be End to end NVMe. <b>Justification</b> : We request for modification to include optimal technologies to deliver the required performance .	Please see corrigendum.
109		we request GIL to add this much required clause to have right storage architecture for this project.	The storage provided must be scale-out architecture with linear performance increase as more capacity & controller/nodes are added in future with single filesystem and namespace for ease of management and scalability .	No change, As per RFP.
110	10. Technical Specification: Nodes for AI training / 10/15	IOPS: Min. 10,00,000 Throughput: Min. 100 Gbps Bidirectional	<b>Modification</b> : IOPS Min 8,00,000. <b>Clarification</b> : Please confirm the asked throughput is Giga bits and NOT Giga Bytes!! <b>We recommend to have at least 20GBps performance for 8,00,000 IOPS. With linear scalability with addition of Capacity and controllers/nodes in future.</b>	Please see corrigendum.
111	10. Technical Specification: Nodes for AI training / 10/16	<b>System Network</b> : InfiniBand for storage delivery	<b>Clarification</b> : Please remove this clause as Leading storage OEM do not have InfiniBand as front end connectivity. We recommend either to use FC @32Gbps for connectivity that will need SAN Switches as well, OR <b>Ethernet 100GBps or more ports for seamless integration with rest of the infrastructure with no bottlenecks.</b>	Please see corrigendum.
112	11.Inference Node, Page- 11/15	The bidder should configured Head/ Master/ Management Node in (1+1) HA mode deliver the solution.	<b>Clarification</b> : Kindly confirm if the Training and Inference nodes will be in separate clusters or a single cluster with Master/Head Nodes in High Availability. Also for HA it is highly recommended to have a separate dedicated shared storage for storing cluster database across the 2 head/Master nodes. Hence we request that a discrete storage array with capacity of about 20TB in raid5 on SSD be asked.	Bidder to propose the best architecture based on the AI workload requirements and ensuring there is no performance loss due to network design. No change
113	11.Inference Node, Page- 11/16	2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator	<b>Modification</b> : 2 x Accelerators per node, each with minimum. 94GB GPU per Accelerator <b>Justification</b> : Since there ask is for Multi Instance GPU, which is not supported on 40GB GPUs , hence request GIL to amend the clause as required for the solution.	Please see corrigendum.
114	Eligibility Criteria: 4   Page No: 1	The bidder should have experience of set up of similar GPU base solution with minimum 75 GPU/CPU based solution with minimum 2000 cores in last 3 Years in last 3 Years in India with below criteria as on last date of bid submission: <input checked="" type="checkbox"/> One project having similar works total value of INR 32 Cr or <input checked="" type="checkbox"/> Two project having similar works total value of INR 20 Cr or <input checked="" type="checkbox"/> Three project having similar works to-tal value of INR 12 Cr Note: Similar works means SITC OF GPU AC-CELERATED with multiple GPU Node.	As this is GPU dominant bid. The requirement should focus specifically on GPU related experience not just CPU bases systems. Request to remove CPU pls.	Please see corrigendum.
115	8. IMPLEMENTATION TIMELINES & PENALTIES   Page No: 7	Supply of the Hard-ware including Li-censes and OEM Warranty Certificate.  T1=T+60 days from the date of issuance of contract over GEM	The delivery timeline should be extended from 60 days to 180 days due to the current market trend and high lead times for GPU systems.	Please see corrigendum. T1=T+45 days from the date of issuance of contract over GEM
116	General	Bid to RA	We request you to remove the RA for this bid. The current bid pertains to highly specialized GPU-accelerated systems. Such solutions are not standardized or commoditized	As per RFP.

## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

117	1 Eligibility Criteria	The bidder should have experience of set up of similar GPU base solution with minimum 75 GPU/CPU based solution with minimum 2000 cores in last 3 Years in last 3 Years in India with below criteria as on last date of bid submission: •One project having similar works total value of INR 32 Cr or •Two project having similar works total value of INR 20 Cr or •Three project having similar works total value of INR 12 Cr Note: Similar works means SITC OF GPU AC-CELERATED with multiple GPU Node.	The bidder should have experience of set up of solution with minimum 1000 CPU cores in last 7 Years in India with below criteria as on last date of bid submission: •One project having similar works total value of INR 32 Cr or •Two project having similar works total value of INR 20 Cr or •Three project having similar works total value of INR 12 Cr	Please see corrigendum.
118	External Storage	The solution should be delivered with 1PB usable post RAID 6 or better configuration, expandable up to 2 PB. The proposed storage array should be configured in no single point of failure including controller (at least 2), cache, power supply, cooling fans, etc. Disks: 1) NVMe SSDs <del>and 2) NL-SAS disks</del> . The storage should be distributed with namespace consistent across nodes. IOPS: Min. 10,00,000 Cache: 1 TB across controllers Throughput: Min. 100 Gbps Bidirectional Connectivity: 8 x 100G ports , 8 x 25G, 8 x 32 G Protocols: NVMe-oF, NVMe/TCP, iSCSI, pNFS, NFS (NFSv3, NFSv4, NFSv4.1 ), CIFS/SMB, S3 Other Features: Ransomware protection , Data at rest and in flight encryption , inline and post process compression and deduplication		Please see corrigendum.
119	3. Scope of Work: Point No 8, Page no 4	The bidder will have to supply Server Rack along with provision of IPDU, TOR Switch, patch panel, cables, SFP modules, any other active/passive components etc. to host the HPC Clus-ter with GPUs at GSDC. Any other component required for the solution proposed by the supplier has to be incorporated for completion of the Solution.	We assume that Electrical & Civil work till Rack IPDU connectivity will be provided by the GIL. Kindly confirm.	Bidder has to provide end to end power and connectivity till rack IPDU for complete solution.
120	General	Management Switch	Management switch is not considered for management of components. Please include 1 No of 24P or 48P Management switch.	Bidder has to consider <b>all required</b> switch for completion of proposed solution. Please see corrigendum.
121	3. Scope of Work: Point No 8, Page no 3	TOR Switch	What is the uplink bandwidth GIL is looking for?. Kindly specify.	please see corrigendum.
122	Nodes for AI training: Total Qty-4 sets. Software Support Page no 10	All necessary and required software, SDK, libraries, tools to cater and run the AI/ML workload should be provide from day 1.	Kindly specify the applications workload and types of AI/ML applications to be hosted.	It will be as per departments requirement.
123	10. Technical Specification: Pg 9	Master Node:	Qty not mentioned, We assume that one Master Node is required. Kindly confirm ?	please see corrigendum.
124	Nodes for AI training: Total Qty-4 sets. Pg 10	Network, Minimum 4 nos. of InfiniBand NDR ports for internode communication, 1 nos. of port for BMC (dedicated LAN port), minimum 1 no. of 1 GbE & 10 GbE (Fiber) port. Additional InfiniBand HDR/NDR200 single/2-port HCA for storage connectivity.	Minimum 2 no. of 1 GbE & 10 GbE (Fiber) ports should be asked for HA.	please see corrigendum.
125	Nodes for AI training: Total Qty-4 sets. Pg 10	<b>Internal Storage:</b> • For Operating System: minimum 1.92 TB M.2 NVMe drives	Min two nos of NVMe Drives should be asked for HA	please see corrigendum.
126	Nodes for AI training: Total Qty-4 sets. Pg 10	<b>HBA Card</b> 16 Gbps Host Bus Adaptor for connecting with storage.	We assume that SAN Switch (having ports of 16Gbps) will be provided by GIL, Kindly confirm. Also, kindly share the make & model of the SAN Switch. Also min two nos of 16 Gbps HBA should be asked for HA	Bidder has to consider SAN switch in their scope for connecting existing Netapp and Hitachi storage. Please see corrigendum.
127	Nodes for AI training: Total Qty-4 sets. Pg 10	<b>External Storage:</b> The solution should be delivered with 1PB usable post RAID 6 or better configuration, expandable up to 2 PB. The proposed storage array should be configured in no single point of failure including controller (at least 2), cache, power supply, cooling fans, etc. Disks: 1) NVMe SSDs and 2) NL-SAS disks. The storage should be distributed with namespace consistent across nodes. IOPS: Min. 10,00,000 Throughput: Min. 100 Gbps Bidirectional	Kindly share the bifurcation of the Storage in terms of % of disk like NVMe. NL-SAS. Kindly share the specification of the external storage in terms of controllers and ports.	please see corrigendum.
128	General	InfiniBand	Kindly specify the quantity of InfiniBand switch and no of ports per switch required.	please see corrigendum.

## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

129	Inference Node: Total Qty-12 sets. Pg 11	<b>Internal Storage:</b> • For Operating System: minimum 1.92 TB M.2 NVMe drives	Min two nos of NVMe Drives should be asked for HA	please see corrigendum.
130	Inference Node: Total Qty-12 sets. Pg 11	<b>HBA Card</b> 16 Gbps Host Bus Adaptor for connecting with storage.	We assume that SAN Switch (having ports of 16Gbps) will be provided by GIL, Kindly confirm. Also kindly share the make & model of the SAN Switch. Also min two nos of 16 Gbps HBA should be asked for HA	Bidder has to consider SAN switch in their scope for connecting existing Netapp and Hitachi storage.  Please see corrigendum.
131	Inference Node: Total Qty-12 sets. Pg 12	<b>Networking Switch</b> Required 10/25/40/100 Gigabit Ethernet switches for complete solution.	Kindly specify the Ports Qty required per switch	please see corrigendum.
132	c. Manpower related SLA and Penalties: P8	The agency has to implement the attendance system and share the attendance report of each person deployed as part of team on monthly basis with the GSDC.	agency meaning bidder ? Bidder needs to install the attendance system	Bidder has to maintain their monthly basis attendance register and report for each deployed person at GSDC.
133	Pg. no 6. point 9 b	In case of failure of proposed solution and non-main taining targeted value, 0.5% of yearly CAMC payment for every hourly. delay in resolution; with max cap of 10 % of total 5 years CAMC value; it will be adjusted from Yearly CAMC payment .	Please revise the clause - "Uptime loss Penalties are proposed on annual contract value"	As per RFP
134	Scope of work, Pt 12, Page 3	If, during the warranty period, any system as a whole or any subsystem has any failure on two or more occasions in a period of 3 months, it shall be replaced by equivalent new equip-ment by the bidder at no cost. All defective items should be replaced / repaired within the SLA period.	The clause should have been that if system or sub-system has REPEATED failures of SAME type for 3 or more occasions in period of 3 months because at times, the identification of failure is complex and 1st or 2nd resolution may still not solve the problem. If system is giving same type of problem repeatedly then its chronic failure and should be replaced with another system. In fact Section 9 at page 7 mentions 3 failures	Revised as , if, during the warranty period, any system as a whole or any subsystem has any failure on three or more occasions in a period of 3 months, it shall be replaced by equivalent new equip-ment by the bidder at no cost. All defective items should be replaced / repaired within the SLA period.
135	Implementation timelines, Pt 3, Page 7	T2 = T1 + 30 days	HLD, LLD may need deliberations hence 30 days will be too aggressive timeline. Request to keep 60 days	As per RFP
136	SLA for uptime, Page 8	SLA will be calculating on monthly/quarterly basis, However, Final penalty deduction on the yearly payment of CAMC value i.e., (4* 3 quarter or 12*1 monthly SLA report penalty will be applied during CAMC yearly payment.)	SLAs to be either calculated monthly or quarterly, hence should be clarified. Also, it would be better if the formula of penalty calculation is provided, especially for uptime related penalties	As per RFP
137	Manpower related SLAs, Page 8	Replacement of a profile by the agency (only one replacement per technical profile – with equal or higher qualification and experience – would be permitted per year)	We assume that such replacements don't include cases where employee leaves on his own as that can't be stopped as per law	As per RFP
138	Manpower related SLAs, Page 8	Replacement of resources by the agency on formal submission of resignation by the resource in the company.	Since overlap of 15 days can't always be guaranteed hence the clause should be - There should be minimum of 15 days of overlap between new deployed resource and replaced resource OR The backup resource provisioned to take care of leave of deployes resource should work as a transit resource between deployed and replaced resource with overlap of 15 days with deployed and 15 days with replaced resource	As per RFP
139	Additional clause	Limitation of liability	Successful Bidder's cumulative liability for its obligations under the contract shall not exceed the value of the charges payable by GIL within the remaining duration of the contract term from the day claim is raised and selected agency shall not be liable for incidental, consequential, or indirect damages including loss of profit or saving	As per RFP
140	Internal Storage Master Node:	• For Operating System: Two numbers of 3.84 TB capacity hot swap Enterprise NVMe SSDs with MTBF of 2 million hours or more, Configure in RAID 1. • For Data: Six numbers of 3.84 TB capacity interface hot swap Enterprise NVMe SSDs with MTBF of 2 million hours or more, configured in RAID 5.	please revise this clause as below: "• For Operating System: Two numbers of 3.84 TB capacity hot swap Enterprise NVMe SSDs, Configure in RAID 1. • For Data: Six numbers of 3.84 TB capacity interface hot swap Enterprise NVMe SSDs configured in RAID 5.  MTBF is OEM specific so please generalise the clause	please see corrigendum.
141	Storage Controller Master Node:	Hardware RAID 0,1, 5, 6, 10, 50, 60 with 8GB cache Flash based cache protection module should be included, should support Gen4.0 PCIe NVMe.	please revise this clause as below: "Hardware RAID 0,1, 5, 6, 10, 50, 60 with min 4GB cache Flash based cache protection module should be included, should support Gen4.0 PCIe NVMe."	please see corrigendum.
142	Network Master Node:	• 1G-Gigabit or higher Ethernet LAN RJ45 ports • One Single Port 200Gb/s NDR InfiniBand Adapter ConnectX-7 (1x OSFP). • One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6	• 1G-Gigabit or higher Ethernet LAN RJ45 ports • One Single / dual Port 200Gb/s NDR InfiniBand Adapter ConnectX-7 (1x OSFP)/Ethernet adapter • One Dual Port 200Gb/s HDR InfiniBand Adapter ConnectX-6/Ethernet adapter	please see corrigendum.
143	Number of GPUs and GPU Communication Nodes for AI training:	8 x Accelerators per node, each with minimum. 80GB GPU per Accelerator. Minimum 800GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL /ROCm with Accelerators,	Please revise as below for wider OEM participation: "8 x Accelerators per node, each with minimum. 80GB GPU per Accelerator. Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL /ROCm with Accelerators,"	please see corrigendum.
144	Internal Storage Nodes for AI training:	• For Operating System: minimum 1.92 TB M.2 NVMe drives • For Data: Minimum 4 * 3.84 TB U.2 NVMe drives	Please revise as below for wider OEM participation: "• For Operating System: minimum 1 TB M.2 NVMe drives • For Data: Minimum 16* 1.9 TB NVMe drives"	please see corrigendum.

## Responses of Pre-Bid Queries (GeM Bid NO. GEM/2025/B/5824041 dated 16.01.2025)

## Bid for Procurement of GPU Compute solution for Gujarat State Data Center, Gandhinagar

145	HBA Card Nodes for AI training;	16 Gbps Host Bus Adaptor for connecting with storage.	Please remove this clause as normally DPUs are required for storage traffic communication	please see corrigendum.
146	System Network Nodes for AI training;	Following N/W are required: 1. Infiniband for compute communication 2. Infiniband for storage delivery 3. Ethernet (10GbE) for cluster orchestration 4. Ethernet (10GbE) for perimeter connectivity 5. Ethernet (1GbE) for in-band management	Please revise as below: "Following N/W are required: 1. Infiniband / Ethernet for compute communication 2. Infiniband / Ethernet for storage delivery 3. Ethernet (10GbE) for cluster orchestration 4. Ethernet (10GbE) for perimeter connectivity 5. Ethernet (1GbE) for in-band management"	please see corrigendum.
147	Number of GPUs and GPU Communication Inference Nodes	2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator. Minimum 600GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators.	please revise as below: "2 x Accelerators per node, each with minimum. 40GB GPU per Accelerator. Minimum 64 GB/s bidirectional communication bandwidth per GPU. Should support Tensor core/Matrix core, CUDA / Stream Processors/ openCL / ROCm with Accelerators."	please see corrigendum.
148	Internal Storage Inference Nodes	<ul style="list-style-type: none"> <li>For Operating System: minimum 1.92 TB M.2 NVMe drives</li> <li>For Data: Minimum 4 * 3.84 TB U.2 NVMe drives</li> </ul>	please revise and generalise as below: "• For Operating System: minimum 1.92 TB NVMe drives/2*960GB M.2 NVMe drives • For Data: Minimum 4 * 3.84 TB U.2/U.3 NVMe drives"	please see corrigendum.