



PRODUCTS		NVIDIA A100 PCIe	NVIDIA A30	NVIDIA A2	NVIDIA A40	NVIDIA A16	NVIDIA A100X	NVIDIA A30X
PNY PART NUMBER		NVA100TCGPU80-KIT	NVA30TCGPU-KIT	NVA2TCGPU-KIT	NVA40TCGPU-KIT	NVA16TCGPU-KIT	NVA100XTCGPUCA-KIT	NVA30XTCGPUCA-KIT
WORKLOAD	DESCRIPTION	Highest Performance Compute	Mainstream Compute	Entry-Level Compact AI	Highest Performance Graphics	Optimized for VDI	Highest Performance Converged Accelerator	Mainstream Converged Accelerator
Recommended Number of GPUs or Converged Boards per Server								
<b>Deep Learning (DL) Training and Data Analytics</b>	For the absolute fastest model training and analytics	<b>4-8 GPUs</b> 80GB: Bn+ parameter models (DLRM, GPT-3)					<b>1-2 cards</b> for multi-node training	
<b>DL Inference</b>	For batch and real-time inference	<b>1-2 GPUs</b> w/ multi-instance GPU (MIG) 80GB: large batch size constrained models (RNN-T)	<b>2-4 GPUs</b> with MIG	<b>1-4 GPUs</b>				
<b>High-Performance Computing (HPC) / AI</b>	For Higher Education Research and scientific computing centers	<b>2-4 GPUs</b> with MIG	<b>2-4 GPUs</b> with MIG				<b>1-2 cards</b> for multi-node workloads	
<b>Render Farms</b>	For batch and real-time rendering				<b>4-8 GPUs</b>			
<b>Graphics</b>	For the best graphics performance on professional VDI			<b>1-4 GPUs</b> for entry-level virtual workstations*	<b>2-4 GPUs</b> for midrange to high-end virtual workstations*		<b>2-4 GPUs</b> for highest virtual desktop and workstation user density**	
<b>Cloud Gaming</b>	For 4K resolution / Android			<b>1-4 GPUs</b> for mobile android	<b>4-8 GPUs</b> (4K resolution)			
<b>Enterprise Acceleration</b>	For mixed workloads, including graphics, ML, DL, analytics, training, and inference	<b>1-2 GPUs</b> with MIG for compute workloads	<b>1-2 GPUs</b> with MIG for compute workloads	<b>1-4 GPUs</b> for balanced workloads*	<b>1-2 GPUs</b> for graphics-intensive workloads*			<b>1 card</b> for compute acceleration with software-defined infrastructure
<b>Edge Acceleration</b>	For differing use cases and deployment locations	<b>1-2 GPUs</b> with MIG	<b>1-2 GPUs</b> with MIG	<b>1-4 GPUs</b> for inference and video workloads	<b>1-4 GPUs</b> for graphics-intensive workloads & AR / VR*		<b>1 card</b> for AI-on-5G with heavy workloads	<b>1 card</b> for AI-on-5G with average workloads
<b>5G vRAN</b>	For low-latency GPU-network communication							<b>1-2 cards</b>
<b>AI-Based Security</b>	For GPU-powered network processing							<b>1 card</b>

\* NVIDIA RTX Virtual Workstation (vWS) software license required for virtual workstation workloads.

\*\* NVIDIA Virtual PC (vPC) software license required for VDI workloads.

PNY Technologies, Inc. 100 Jefferson Road, Parsippany, NJ 07054 | Tel 973-515-9700 | Fax 973-560-5590 | [www.PNY.com](http://www.PNY.com)

Features and specifications subject to change without notice. The PNY logo is a registered trademark of PNY Technologies, Inc. All other trademarks are the property of their respective owners.

©2021 PNY Technologies, Inc. All rights reserved.

