# HPE ProLiant Compute Gen12 achieves multiple world records in AI inference benchmarks

Achieving #1 rankings across more than 23 MLPerf™ Benchmark models for HPE ProLiant Compute DL384 Gen12 and HPE ProLiant Compute DL380a Gen12 servers

## HPE ProLiant Compute DL384 Gen12 delivered 13 new world records on MLPerf Inference: Datacenter v5.0[1]

**#1**
**Stable Diffusion XL (SDXL)**
Advanced image generation model that produces high-quality, detailed images from text descriptions

**#1**
**Llama2-70B-99**
High-precision variant of Llama 2-70B, optimized for enhanced accuracy and reliability in AI reasoning tasks

**#1**
**Llama2-70B-99.9**
Ultra-high-precision version of Llama 2-70B, designed for near-perfect consistency and minimal errors in complex AI applications

**#1**
**Mixtral-8x7B**
High-quality sparse mixture of experts model known for its exceptional performance and fast inference speeds

## MLPerf Inference: Datacenter v5.0 benchmark results on HPE ProLiant Compute DL384 Gen12 server[1]

**#1**

Superior AI power: High performance server with the NVIDIA® GH200 NVL2 accelerator[2]

| Benchmark tests | Server[3] | Offline[4] |
|---|---|---|
| Stable Diffusion XL (SDXL) | 4.34 | 5.02 |
| Llama2-70b-99 | 8674.58 | 9362.85 |
| Llama2-70b-99.9 | 8674.58 | 9362.85 |
| Mixtral-8x7B | 15570.80 | 16703.40 |

In a groundbreaking first for our MLPerf results, the HPE ProLiant Compute DL384 Gen12 server with dual-socket NVIDIA GH200 Grace Hopper™ Superchip 144 GB delivers 2x the performance of our single-socket setup with linear scaling.

HPE ProLiant Compute DL384 Gen12 server delivers low-latency datacenter inference with scalable 1P and 2P systems, nearly doubling AI performance in most scenarios.

## HPE ProLiant Compute DL380a Gen12 Server top performer on four benchmarks[1]

**#1**
**Image classification[5]**
ResNet50 Server and Offline benchmarks

**#1**
**Object detection[6]**
Retinanet Server benchmark

**#1**
**LLM summarization[7]**
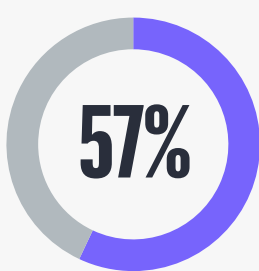GPT-J 99 Server benchmark

**#1**
**LLM summarization[8]**
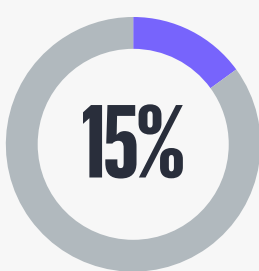GPT-J 99.9 Server benchmark

**#1**

Superior performance over other vendors[1]
HPE ProLiant Compute DL380a Gen12 server

**DLRM-v2-99—Offline[9]**

**Superior performance**

**ResNet50—Server[5]**

**57%**

**15%**

Better than the next top-performing server with 141 GB GPUs

**Visit HPE.com**

**Learn more at**
HPE.com/ProLiant

Chat now