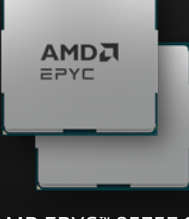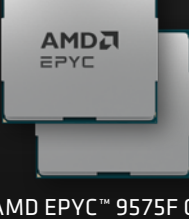**AMD**
together we advance_

# HOW TO GET UP TO
## *20% MORE AI PERFORMANCE*
# FROM HIGH-PERFORMANCE GPUs[1,2]

**5th Generation AMD EPYC™ CPUs** include higher-frequency models designed specifically for hosting accelerator platforms. These CPUs excel at orchestrating data movement and managing multiple virtual machines—critical capabilities that extract more performance from GPU platforms.[1,2]

## The right host makes a difference

### SELECT AMD EPYC HOST NODE CPUs INCREASE INFERENCE AND TRAINING THROUGHPUT

| HOST CPU | GPU PLATFORM | INFERENCE | TRAINING |
|---|---|---|---|
| 2P AMD EPYC™ 9575F CPUs (128 total cores) | 8x AMD Instinct™ MI300X GPUs | Llama 3.1-70B @ FP8 — Up to **8%** more tokens/sec[3] — 700K more tokens/second inference on a 1K node cluster of AMD Instinct GPUs | DeepSpeed 0.14.0 @ FP8 Stable Diffusion XL v2 training set — Up to **20%** more samples/sec[4] |
| 2P AMD EPYC™ 9575F CPUs (128 total cores) | 8x NVIDIA H100 GPUs | Llama 3.1-70B @ FP8 — Up to **20%** more tokens/sec[5] | Llama 3.1-8B @ BF16 Max sequence length 1024 — Up to **15%** more samples/sec[6] |

Performance compared to 2P Intel® Xeon® Platinum 8592+ (128 total cores) hosting the same GPUs running identical workloads.

## Built to boost AI accelerator performance

5 GHz EPYC 9575F max boost is 28% higher than Intel® Xeon® Platinum 8592+.[7]

- 64 energy-efficient cores
- 12 DDR5 memory channels
- Up to 256 MB cache
- Up to 160 PCIe® Gen5 lanes (2P)

## A range of options for hosting GPUs

With frequencies up to 5 GHz and support for up to 6 TB of memory, 5th Generation AMD EPYC CPUs provide multiple models specifically designed for GPU clusters.

| PROCESSOR | CORE COUNT | MAX BOOST FREQUENCY |
|---|---|---|
| 9575F | 64 | 5 GHz |
| 9475F | 48 | 4.8 GHz |
| 9375F | 32 | 4.8 GHz |
| 9275F | 24 | 4.8 GHz |
| 9175F | 16 | 5 GHz |

# 5TH GENERATION AMD EPYC™ CPUs:
## *THE BEST CPU FOR ENTERPRISE AI[8]*

Gain industry-leading server performance for AI, enterprise, and cloud workloads with 5th Generation AMD EPYC processors.

**Explore AMD EPYC**

**AMD**
together we advance_