# Train fast, learn faster

## with the right infrastructure

Staying ahead of the competition requires enterprise AI. And infrastructure that breaks barriers.

Pair the Power AC922 with the IBM Watson Machine Learning Accelerator to help reduce model training times, accelerate iterations and improve insights.
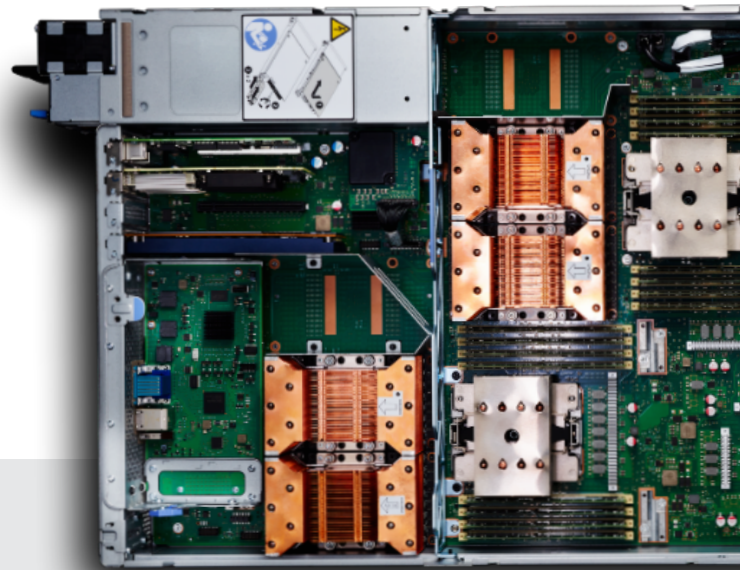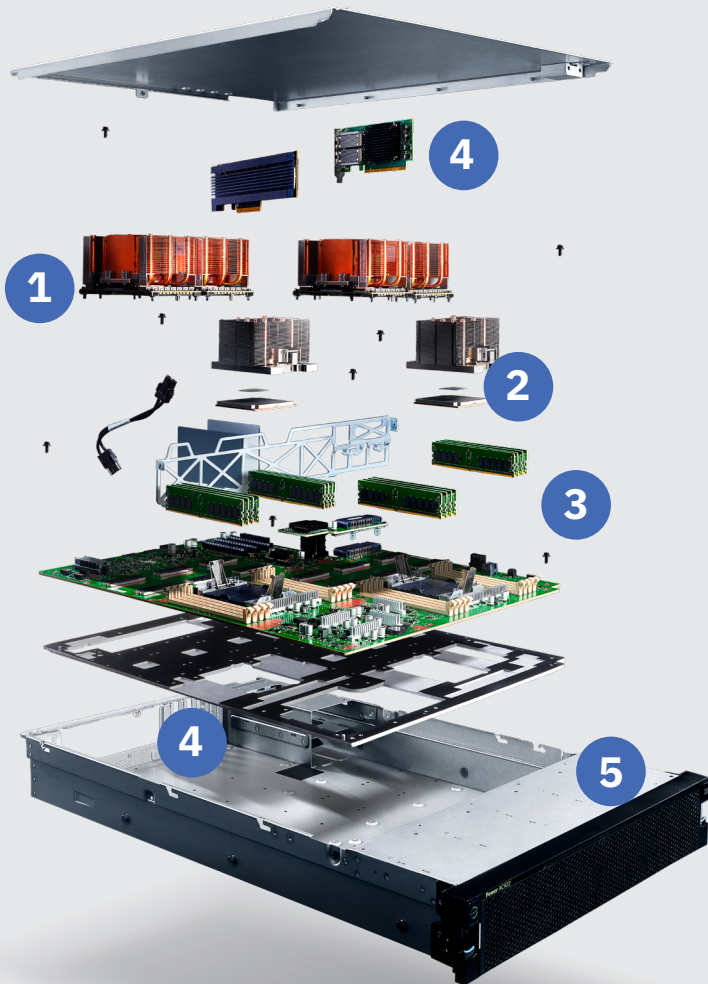
## Solve at the speed of Summit

The IBM® Power System AC922, which powers the world's fastest supercomputer, is purpose-built for AI training.

**Power AC922 + IBM Watson Machine Learning Accelerator**

3.7x **faster training for Caffe[1]**

3.8x **faster training for Chainer[2]**

46x **faster Machine Learning iterations with SnapML[3]**

**IBM POWER9™ + NVIDIA® NVLink™**

5.6x **faster data throughput[4]**

- Up to 6 NVIDIA® Tesla® V100 GPUs + 2 IBM POWER9 processors

- State-of-the-art IO subsystems to handle massive data volume

- NVIDIA NVLink between CPUs and GPUs as well as GPUs

## IBM Power System AC922

**1** GPUs – Up to 6 NVIDIA Tesla V100 GPU processors

**2** CPUs – 2 POWER9 processors with up to 44 cores

**3** System memory – 2 TB max with 16 memory DIMM slots

**4** 4x PCIe Gen 4 slots

**5** Storage – 2 SFF (2.5")
SATA drives, Max 4 TB (HDD)
Max 7.68 TB (SSD)

## Quickly build, train and retrain AI models using a server engineered to be the most powerful training platform

## Get started now

https://www.ibm.com/it-infrastructure/power/enterprise-ai