



GPUUltima

The OSS GPUUltima is a petaflop compute accelerator in a single rack. The OSS GPUUltima supports up to 128 double-wide interconnected accelerators. Consuming only 56 kilowatts of power, the platform uses 94% less power than other petaflop solutions. In addition the single-rack solution occupies 96% less space than other solutions. The OSS GPUUltima employs 64 OSS PCIe Gen3 x16 adapters and cables between the OSS servers and the OSS High Density Compute Accelerators (HDCA) each containing 16 NVIDIA K80 GPU's and 32 Mellanox EDR 100Gbs InfiniBand adapters and cables connecting the GPUs through a 36-port InfiniBand switch. One Stop Systems' PCIe expansion systems provide high speed connectivity between the CPU and IO cards.

Features

- Single 19" rack
- Power input and distribution
- 8 OSS compute accelerators
- 16 dual socket servers
- 128 NVIDIA dual GPU cards
- Infiniband Switch
- Ethernet Switch

Specifications

Rack	<ul style="list-style-type: none"> o 42U tall rack with extra depth o Also available in 24U, 44U and 48U tall versions o Optional doors, etc.
High Density Compute Accelerator	<ul style="list-style-type: none"> o Supports 16 double-wide GPU cards o Four PCIe over cable connections to host servers o 6,000 watts n+1 power o High-speed push-pull fan configuration and channeled airflow
Server	<ul style="list-style-type: none"> o Supermicro 1U Server o Dual socket with Intel Xeon "Broadwell" CPUs o Two 3.5" hot-swap disk drives o Four Gen 3 x16 PCIe expansion slot plus 2 x8 slots o Redundant 2,000 watt power supplies
Infiniband Switch	<ul style="list-style-type: none"> o Mellanox 36 port Infiniband switch o EDR 100Gb/s, QSFP connectors o 1U form factor
Ethernet Switch	<ul style="list-style-type: none"> o Cisco 48 port switch o 1U Form Factor
GPUs	<ul style="list-style-type: none"> o NVIDIA K80 o Contains 2 GPU components o 4,992 cores per card o 8.74 Tflops per card (boost enabled) o 24GB memory o PCIe Gen 3 x16 interface to backplane o 300 watts

PCIe Over Cable Interface Card	<ul style="list-style-type: none"> o PCIe Gen 3 x16 interface o Two cards used per server
Infiniband Interface Card	<ul style="list-style-type: none"> o EDR 100Gb/s, QSFP connectors o Single port version, dual port also available o Two cards used per server
Power Distribution Unit Option 1	<ul style="list-style-type: none"> o Tripp-Lite Monitored PDU o 12.6kW power o Input: 208V 3 phase, 60A o Power monitoring via display and Ethernet o System utilizes 8 PDUs o 100kW total power ~ 80% over-provisioned
Power Distribution Unit Option 2	<ul style="list-style-type: none"> o Tripp-Lite Monitored PDU o 27.6kW power o Input: 380/400V 3 phase, 63A o Power monitoring via display and Ethernet o System utilizes 4 PDUs o 110kW total power ~ 97% over-provisioned
Cables	<ul style="list-style-type: none"> o Power cables – 59 <ul style="list-style-type: none"> ✓ 32 for servers ✓ 24 for compute accelerators ✓ 1 for Ethernet switch ✓ 2 for Infiniband switch o PCIe cables – 32 <ul style="list-style-type: none"> ✓ 32 for servers to compute accelerators ✓ Gen 3 x16 – 1 meter long o Ethernet cables – 48 <ul style="list-style-type: none"> ✓ 32 for servers to Ethernet switch ✓ 8 for compute accelerators to Ethernet switch ✓ 8 for PDUs to Ethernet switch o Infiniband cables – 32 <ul style="list-style-type: none"> ✓ 32 for servers to Infiniband switch
Optional Liquid Cooling	<p>CoolIT Direct Contact Liquid Cooling Model DCLC</p> <p>Custom solution providing special heat sinks with liquid cooling onto each GPU card</p> <p>Tubing to front of each GPU canister</p> <p>Plumbing within rack</p> <p>External pump and heat exchanger</p>

Architecture

