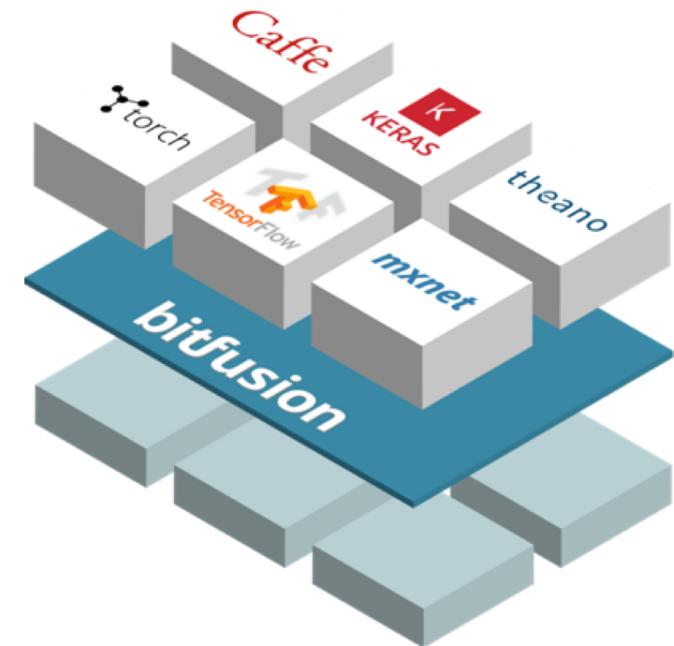


bitfusion

Elastic AI Platform

Jun-2018

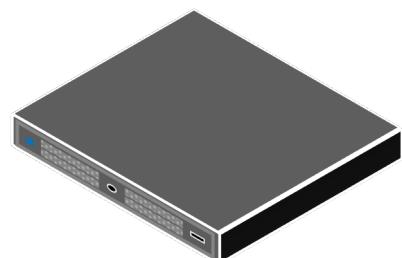




Artificial Intelligence Drives Markets over \$100B

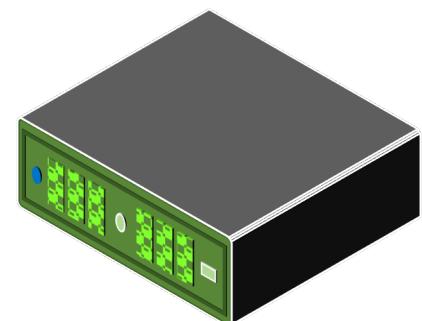


INFRASTRUCTURE



x86 Server

20-50X
Performance Gain



GPU Server



**How many total x86 Servers (Globally)
are deployed:**

- >100M
- 80M – 100M
- 60M – 80M
- 40M – 60M

x86 Server



200 : 1
R A T I O

How many total x86 Servers (Globally)
are deployed:

- >100M
- 80M – 100M
- 60M – 80M
- 40M – 60M

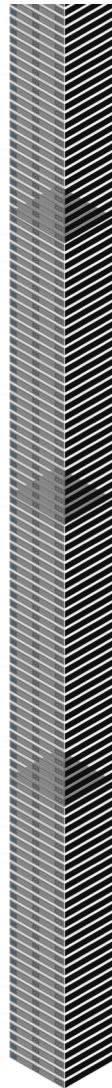
x86 Server

How many total GPU Servers (Globally)
are deployed:

- >1M
- 0.75M – 1M
- 0.5M – 0.75M
- 0.25M – 0.5M



GPU Server



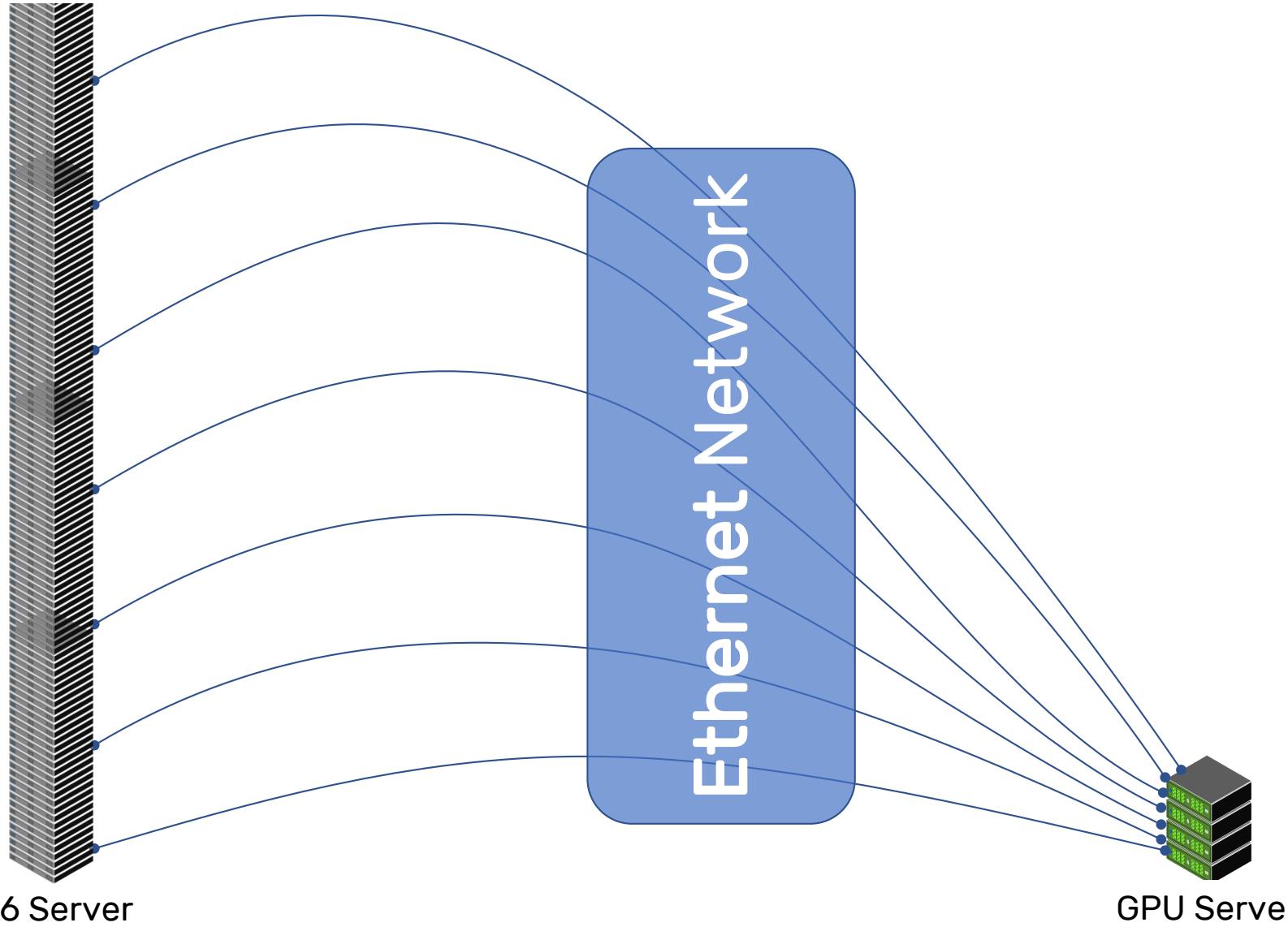
x86 Server

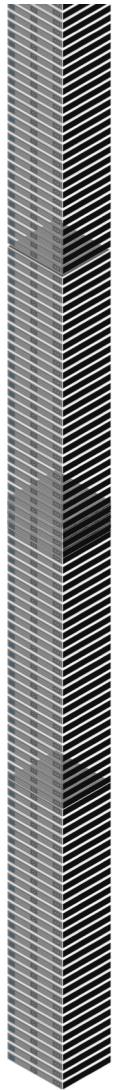
1 Billion man-year Software Eco System
\$7,000 ASP
500w Typical Power

0.05 Million man-year Software Eco System
\$50,000 ASP
3000w Typical Power



GPU Server

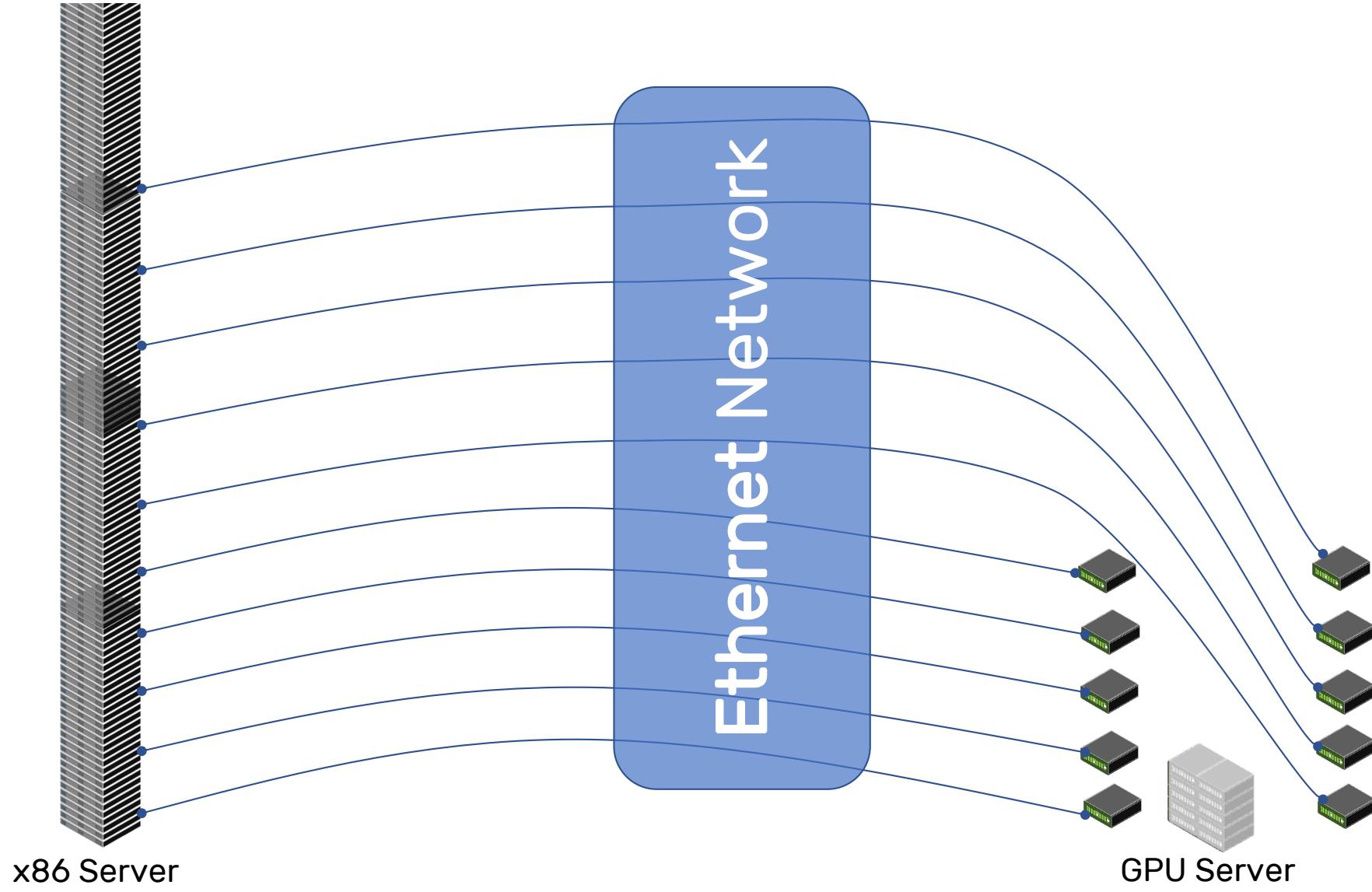


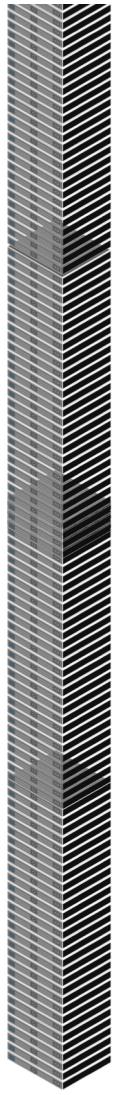


x86 Server

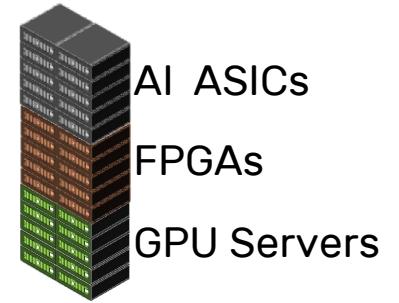


GPU Server





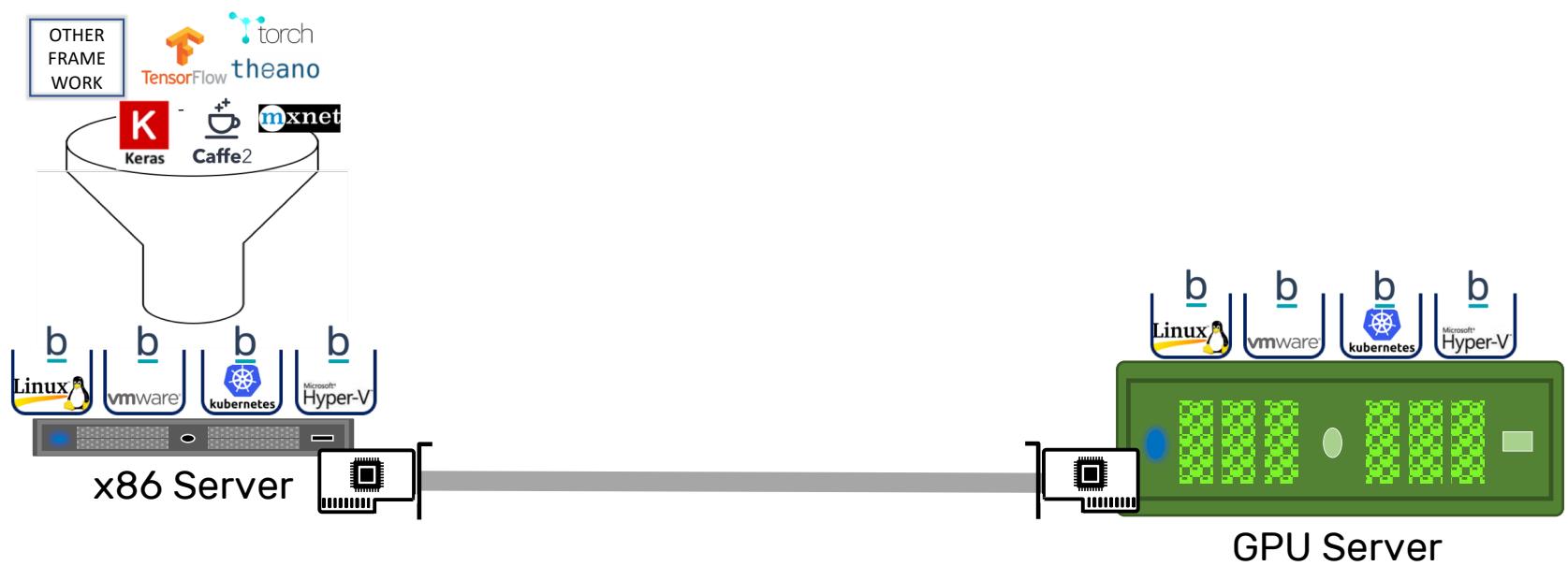
x86 Server

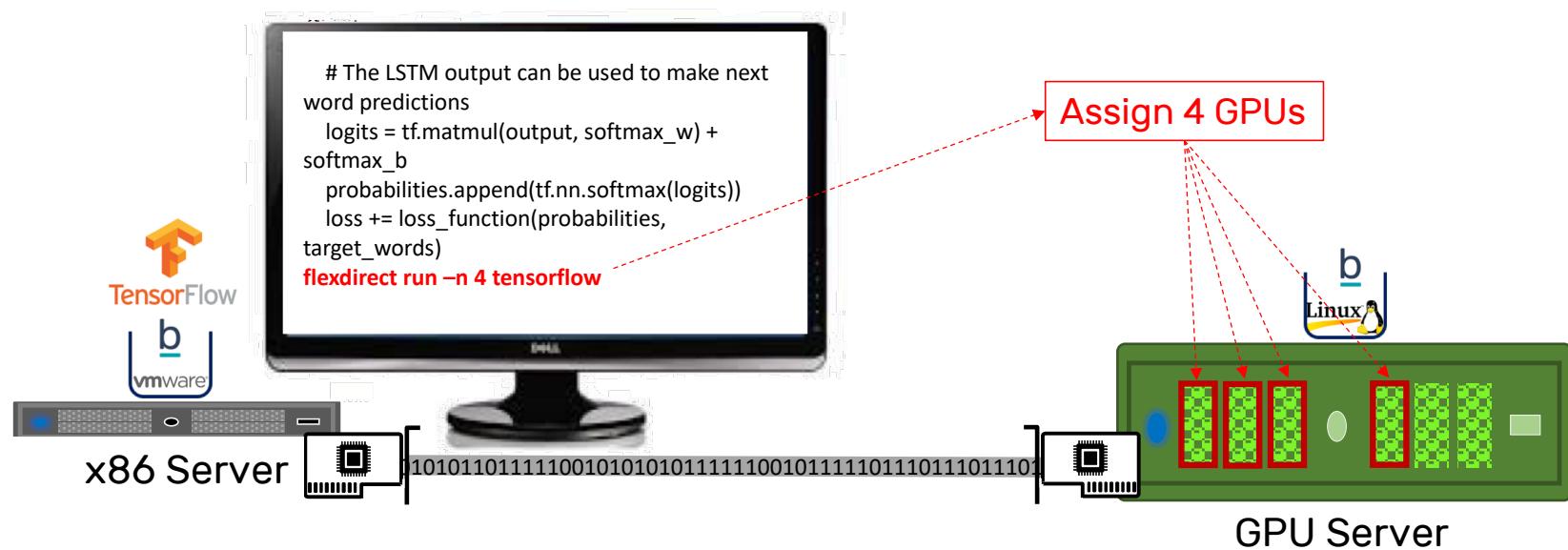


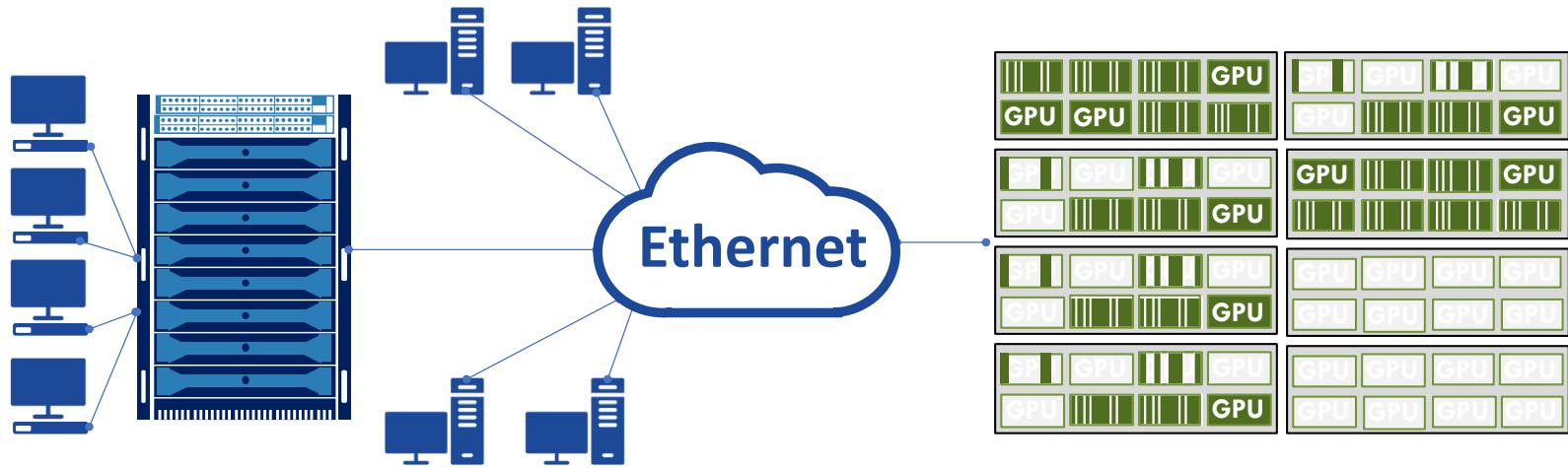
AI ASICs

FPGAs

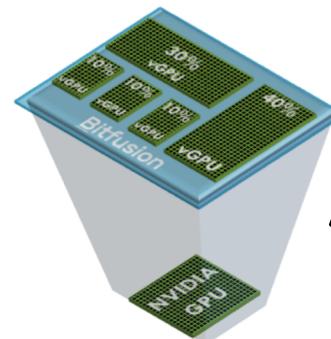
GPU Servers





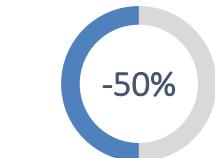


AI Attached Network



AI Hypervisor

Industry's First Elastic AI Infrastructure



Lower Capex



Lower Opex



Productivity



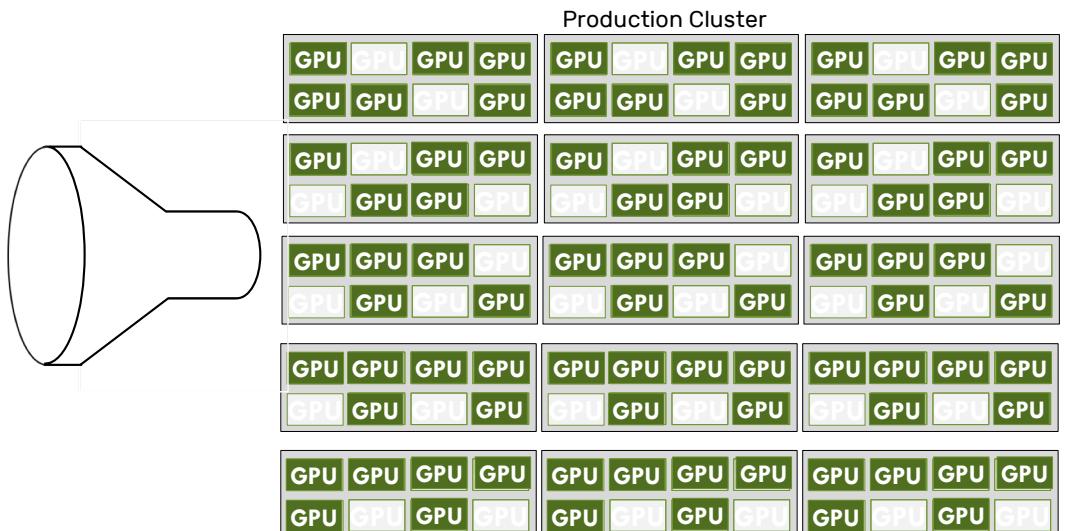
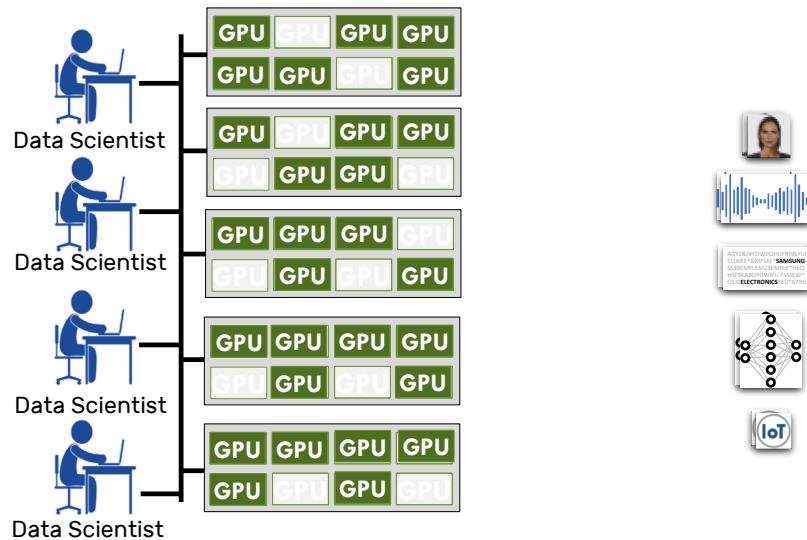
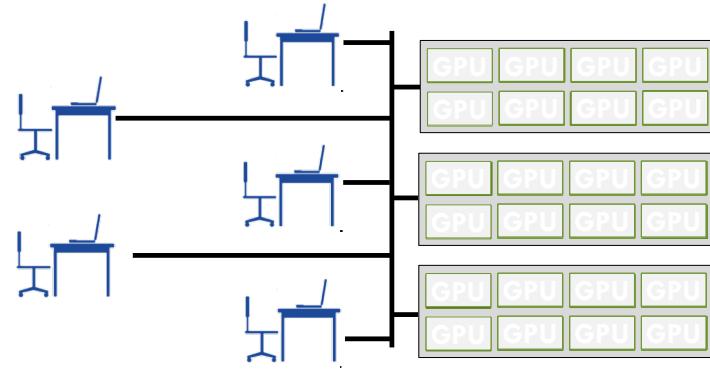
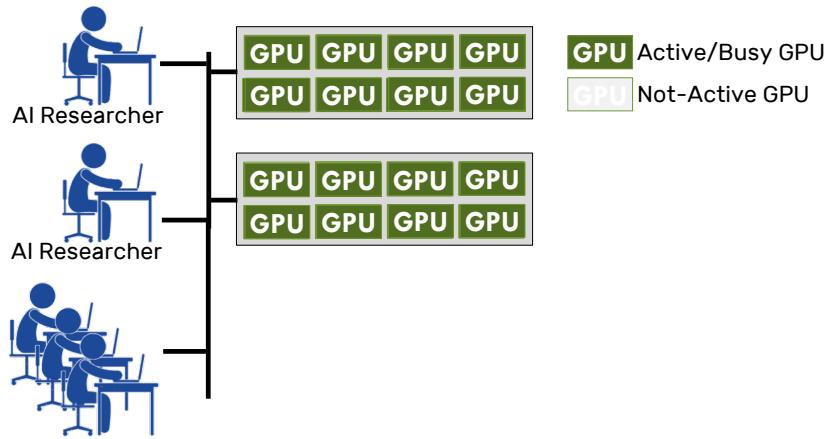
Flexibility



Agility

Virtual Elastic GPU cluster from all the scattered deployed GPU servers in the enterprise. Any user, framework and compute server (without GPU resources) can attach instantaneously to a remote fractional GPU, single GPU or group of GPUs in the virtual cluster, run the AI code, and detach.

Bitfusion virtualize the capacity and the location of GPUs, and make them accessible to any compute machine in the network. The GPUs do not need to be viewed as physical entities, and 1/20, 1/7, 1/4, or any fraction of GPU can be assigned to a workload, on demand.



| | | | |
|-----|-----|-----|-----|
| GPU | GPU | GPU | GPU |
| GPU | GPU | GPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | GPU | GPU | GPU |
| GPU | GPU | GPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | GPU | CPU | CPU |
| GPU | GPU | GPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| GPU | GPU | GPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | CPU |
| GPU | GPU | GPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| GPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| CPU | GPU | GPU | CPU |

| | | | |
|-----|-----|-----|-----|
| GPU | GPU | GPU | CPU |
| CPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | GPU | GPU | GPU |
| CPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | GPU | GPU | GPU |
| GPU | CPU | GPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| GPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| GPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| GPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| GPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| GPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| GPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | CPU |
| CPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | CPU |
| CPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | CPU |
| CPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| CPU | GPU | CPU | GPU |

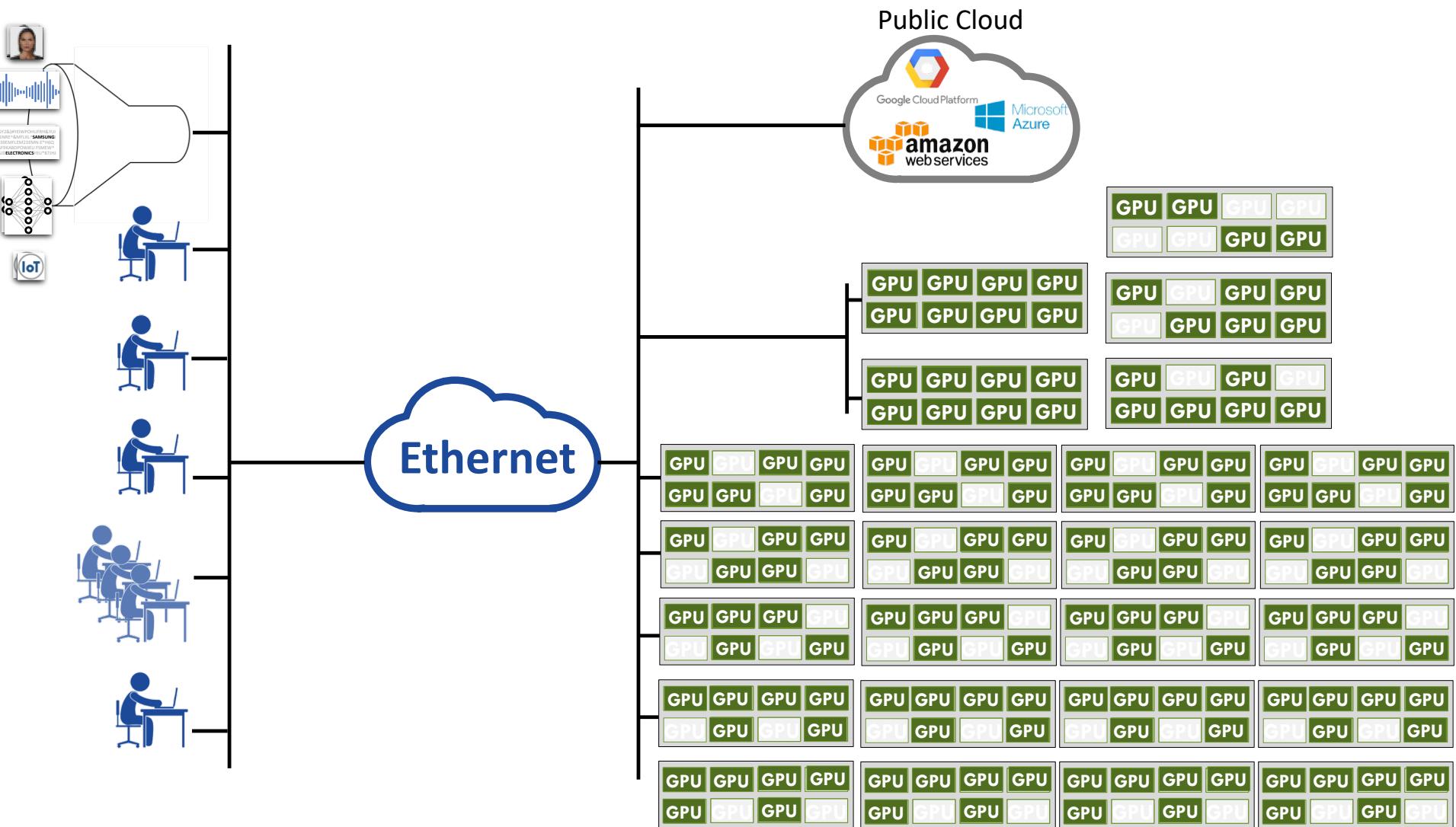
| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| CPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| CPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| CPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| CPU | GPU | CPU | GPU |

| | | | |
|-----|-----|-----|-----|
| GPU | CPU | GPU | GPU |
| CPU | GPU | CPU | GPU |



Any Application/Framework

Any OS/Hypervisor

Any Network

Any Cloud

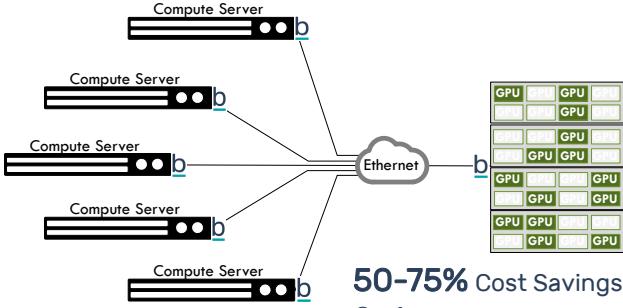
Any Scale

bitfusion - The **Elastic** AI Company



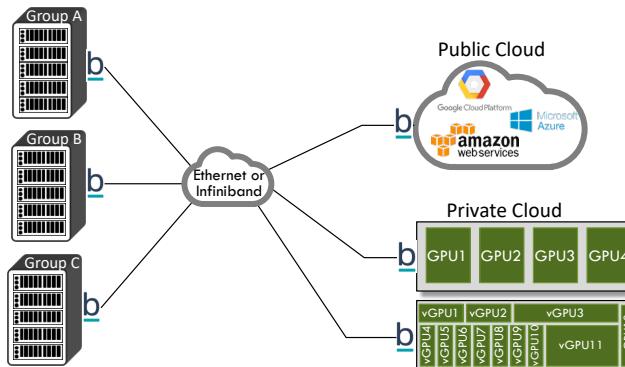
Dave Simonds

Dynamic and Remote Attach



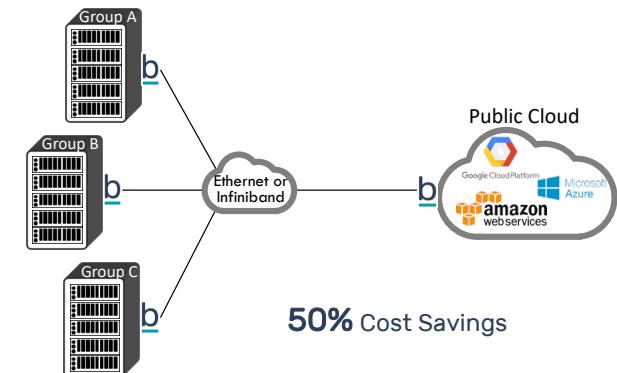
Bitfusion enables compute servers to attach dynamically at application runtime, only when GPU resources are needed to remote GPUs. Multiple Compute servers can attach dynamically to a single or many GPU server. Dynamically attach and detach allows Server consolidation, increase 2x-4x utilization of GPU (or serving 2x-4x more users)

Private & Public Cloud Support



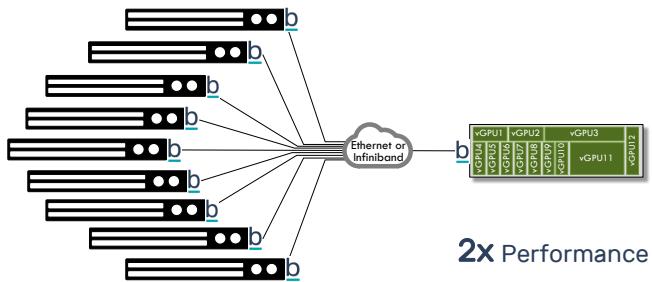
Bitfusion runs in user space, and supports public and private cloud configurations. Many use cases, allow peak GPU demand to burst to the public cloud and extend GPU elasticity in private cloud to public cloud

Optimization of Multi Cloud



Bitfusion runs in user space, and support multi-cloud (public and private). Users can partition GPU server instances to many vGPU and serve more users and applications. Users can pair any compute and GPU instances.

Partial (Virtual) GPUs



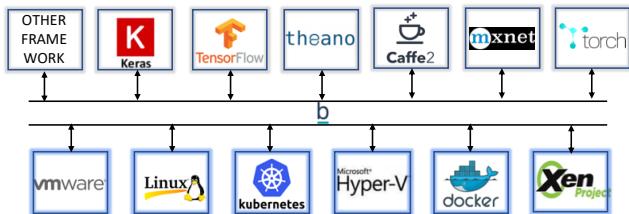
Bitfusion partition physical GPUs to partial virtual (not equal) GPUs. Remote (or local) compute servers can remotely attach to the partial GPUs dynamically. Partial GPUs have demonstrated highest throughput when workload don't consume all GPU memory and/or compute. Many inference applications have a good fit to this topology. Bitfusion can partition a GPU to non-equal parts to allow experimentation and optimization.

GPU Server Consolidation



Bitfusion dynamic attach and detach, and auto-shutdown allows GPU server consolidation, and deliver 4x better utilization (and cost savings). It will translate to similar savings of Opex and Energy savings.

Any Hypervisor, Container or OS



No Management, Provisioning, Training or Administration Costs

Bitfusion is a transparent software layer that supports any upstream AI framework or container environment and any downstream hypervisor, OS or bare-metal. Bitfusion provides flexibility to the software environment while providing the cost, performance and utilization multipliers