# HPC Network Design in Finance

Shawn Hall

# Jump Trading

- Privately-owned proprietary trading firm, established 1999

- World-wide operations
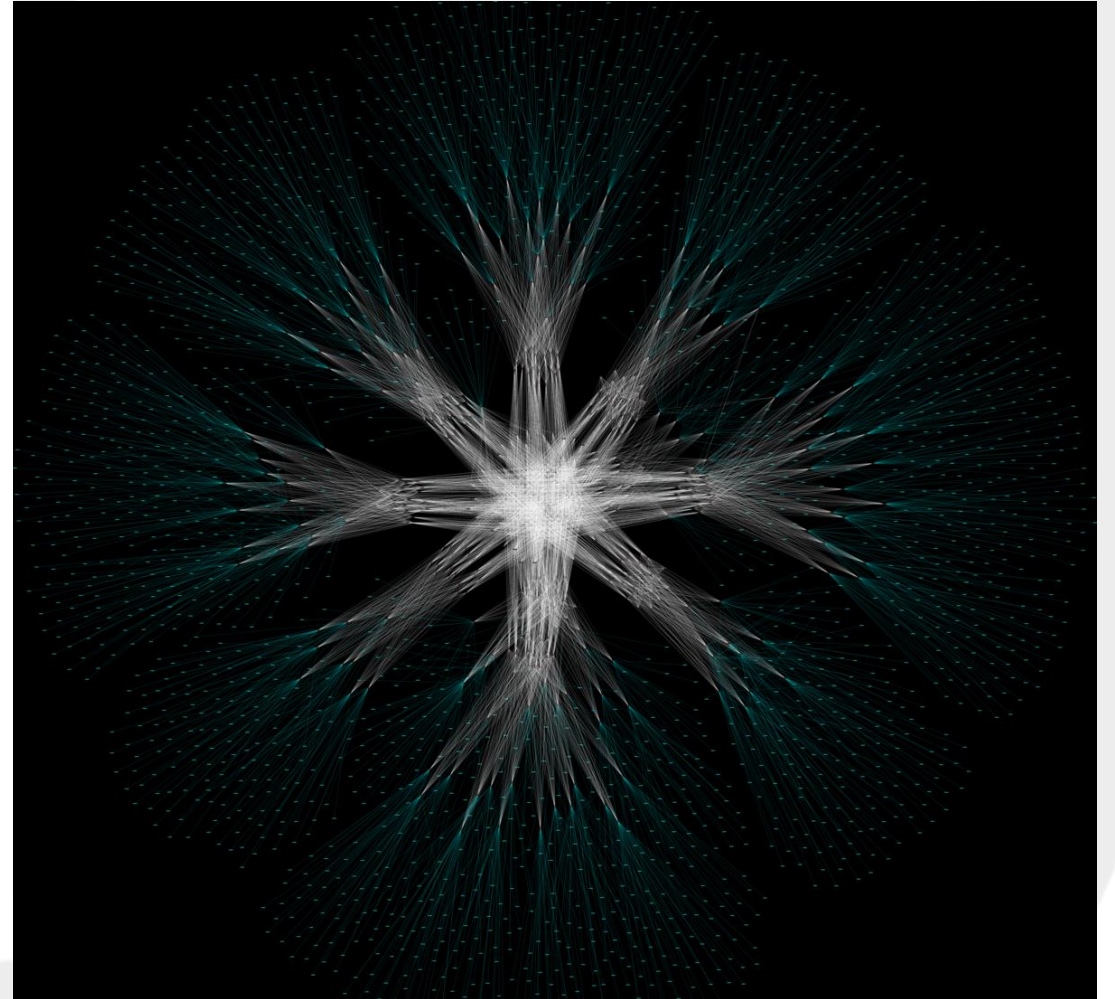  - 12+ offices across US, EU, Asia Pacific

# HPC at Jump

- Research environment with clear correlation to Jump's success

- Platform where we develop and optimize trading strategies

- HPC is critical to operating our business

- Sophisticated data and compute-intensive research workflows

- Technologically competitive with some of the largest publicly known research systems in the world

# Agenda

- **Where we started**

- Where we're at

- Where we're going

# Previous Jump HPC Fabrics

- FDR Infiniband

- Large Clos 5 fabrics

  - Top of rack leafs

  - Director class spines

- Had a terrible time with inter switch links and congestion
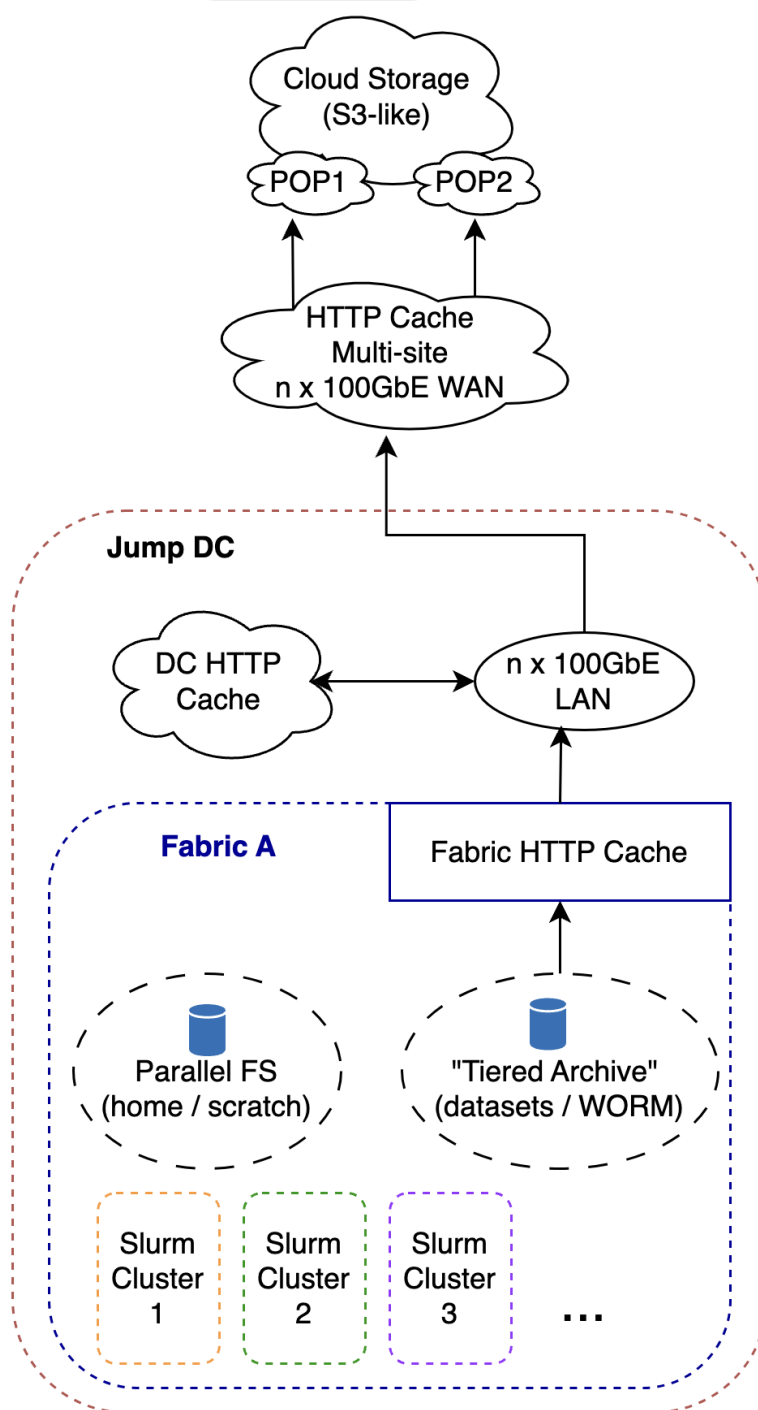
# Agenda

- Where we started

- **Where we're at**

- Where we're going
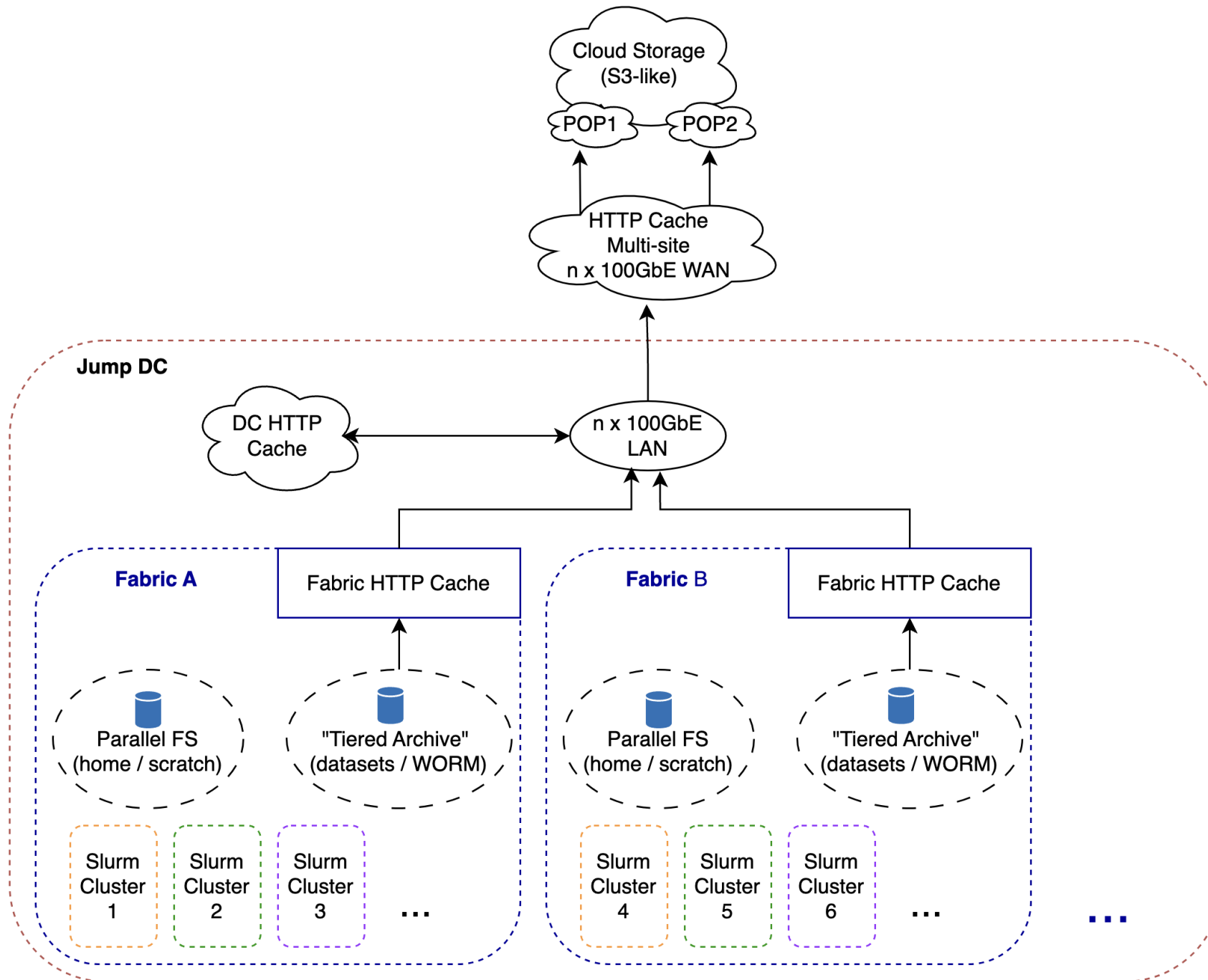
# Current Jump HPC Fabrics



- Textbook HPC components

  - RDMA-capable fabric

  - Parallel filesystem

  - Workload manager (Slurm)

- Add: Global write-once read-many storage system[1]

  - Read-only filesystem presentation (CVMFS)

  - Backed by HTTP caches and cloud storage

  - rsync-like write interface for users
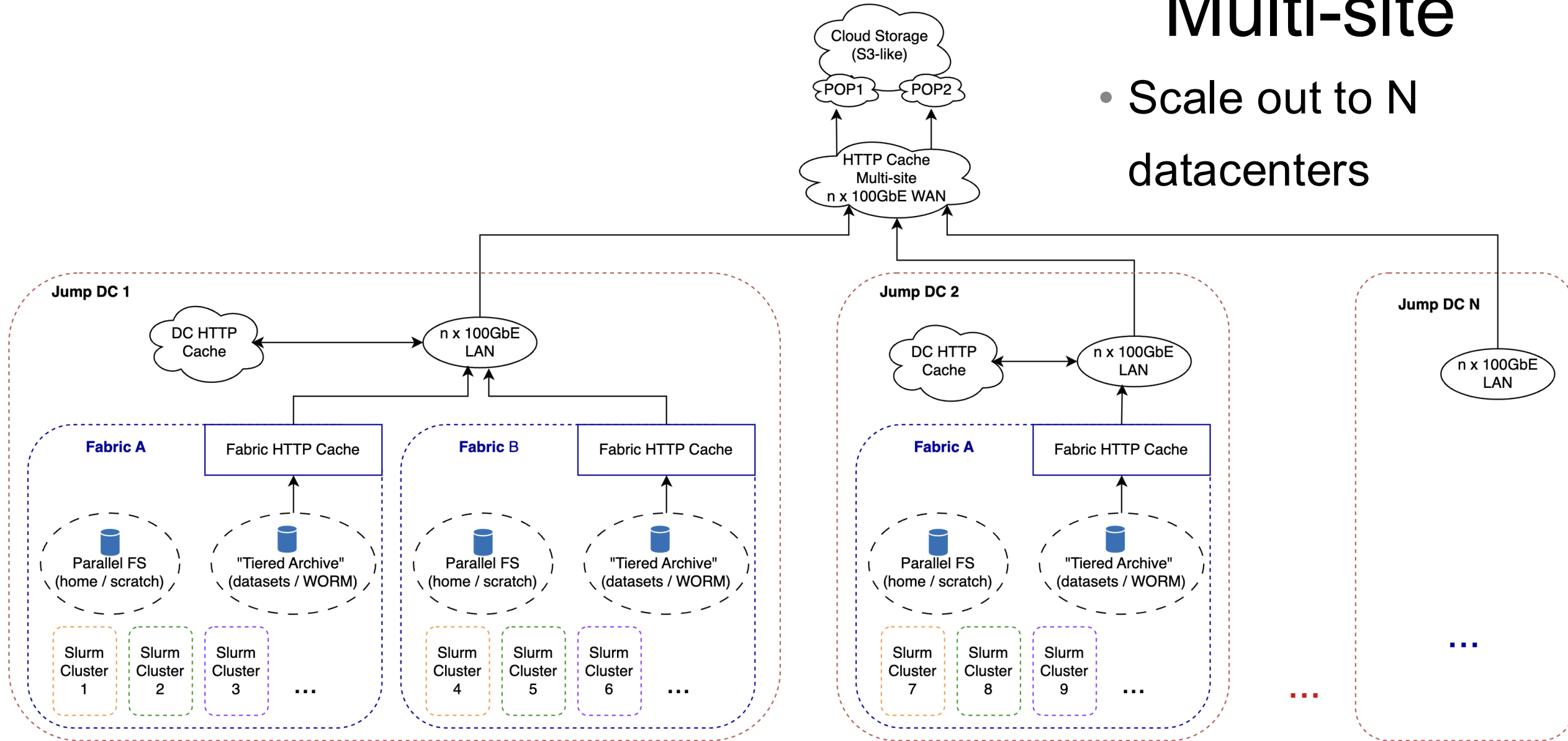
# Multi-fabric



- Scale out to N fabrics in a DC

- Contains blast radius of fabric and parallel FS issues

- Data sharing among fabrics via WORM file system only

# Multi-site

- Scale out to N datacenters

# Some Current Jump HPC Fabrics

- ## CS8500 "Manta Ray"

- ## 800 IB ports
  - ### Splits into 1600 ports at half speed

- ## EOL product – no longer sold, support ends 2029

Port management display

16x Power supplies

I/O Panel

20x Leaf modules

2x Water manifolds

Air ventilation

2x Management modules

20x Spine modules

Cooling Distribution Unit (CDU)

# Somewhat Unique Challenges

# Geographically Distributed HPC

- HPC pods are spread across multiple geographically distant data centers

- Disaster recovery is critical – HPC is core to our business

- More opportunities to find power

- Difficult for data movement
  - WORM file system solves many issues with this

- Lots of network connectivity needed



10 DISASTER SCENARIOS TO CONSIDER

- Natural Disaster (Earthquakes, Floods, Hurricanes, Tornadoes)
- Cyber Attack (Ransomware, Data Breaches, DDoS)
- Fire (Fires in the Office or Data Center)
- Supply Chain Disruption (Supplier/Logistics Issues)
- Terrorist Attack (Bombing, Armed Attacks, Vandalism)
- Power Outage (Extended Power Failures)
- Regulatory Change (Sudden Changes in Laws and Regulations)
- Human Error (Accidental Deletion, Incorrect Configurations)
- Pandemic (Health Crises Like COVID-19)
- Hardware/Software Failures (Critical System Failures)

ITProToday

https://www.itprotoday.com/disaster-recovery/introduction-to-it-disaster-recovery-planning
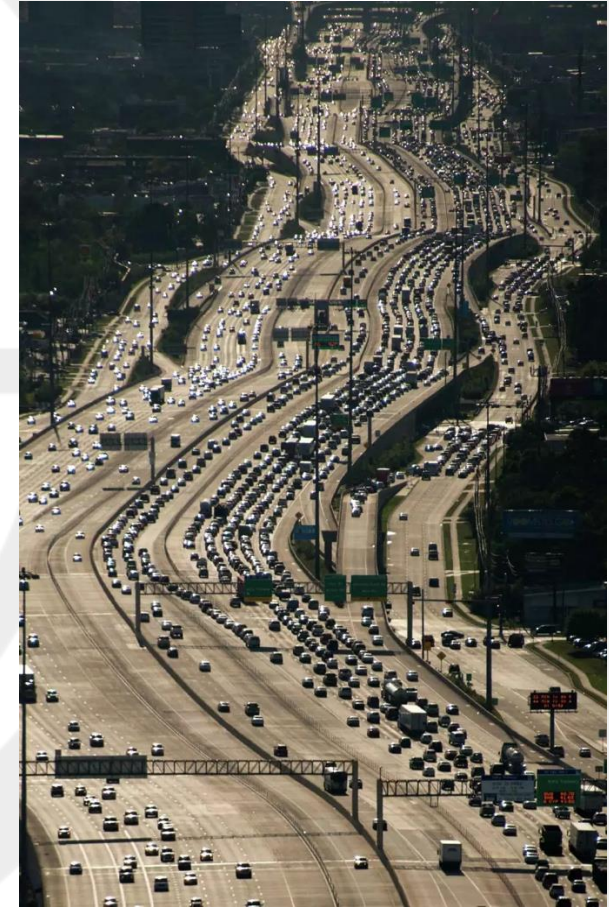
# Dynamic Clusters

- Jump clusters behave like a public cloud

- Many clusters are multi-tenant (different teams/projects)

- Nodes reboot when moving between tenants

- Rolling reboots used to make invasive but regular changes

  - Allows Jump to be nimble and avoid frequent maintenance windows

  - Jump is a fast paced business with quick iteration times – availability of systems and time to resolution are critical for us

- Frequently rebooting compute nodes can cause network issues

# Mixed Use Networks

- IB verbs and IPoIB traffic

- All to all and N to M traffic patterns, and IO elephant flows

- Highly multi-tenant fabric

- Heavy hitter applications segregated via IB virtual lanes

- Reliance on multicast

- Small percentage of traffic generated by MPI libraries

https://www.houstonchronicle.com/politics/texas/politifact/article/World-s-widest-freeway-is-not-where-Turner-thinks-7248455.php

# We Like Big Switches

- FDR experience has made us fearful of leaf spine networks for RDMA

- Flaky link on a network is made exponentially worse by tightly-coupled parallel file systems

- We've had good experiences with Manta Rays
  - *knocks on wood*

- Limited options for director-class switches exist

# Agenda

- Where we started

- Where we're at

- **Where we're going**

# Future Questions

- *Which RDMA network technology?*

- Ethernet
  - Many different flavors exist today

- Infiniband

- Omni-Path

- Others?

# Future Questions

- *Will compute solutions drive network choice/design?*

- Reference architectures guide you toward certain technologies

  - We want to avoid the road less traveled

- Certain compute platforms come integrated with networking

- Should we have separate compute and storage fabrics?

# Future Questions

- *What network architecture to build?*

- Fully independent fabrics like now, but multi-homed storage?

- Islands of compute with limited uplink, multi-homed storage?

- Fully non-blocking cross-pod fabric?

  - How can we get more comfortable with leaf-spine again?

- Does the choice of network technology change the design?

- What can we manage without a hyperscaler-sized network team?

# Future Questions

- *How fast does the storage need to be?*

- Estimate based on network, checkpointing, simple applications, vendor recommendations?

- Can we use node-local NVMe to reduce load on storage system

  - How to best use node-local NVMe?

- POSIX vs. S3?

# Summary

- Jump is a fast paced environment, with short iteration times and a need to evolve quickly

- HPC is a core asset for Jump

- Lots of decisions to make for future networks

- Many strong technology options to pick from

- Lots to learn and test in order to make informed decisions

- Need to mitigate risks and plan for our future

# Jump is Hiring!

- Do these problems sound interesting?

- https://www.jumptrading.com/careers/

# Q&A