

Shift to Server Side Computing

- Computing is shifting to very small user-side and very big devices server-side

- Improved UX
- Advantage for Vendors
 - Simpler to make SW changes and improvement
- Compute-intensive workloads는 server가 처리하기 훨씬 쉽다.

WSC

: Warehouse Scale Computer

- datacenter가 very big data.

Data Center

- **co-located machines** that share power, security, environment, and network connectivity needs
- applications: a few binaries, each running on a dedicated, isolated HW infrastructure
- **heterogeneous hardware and system software**
- **partitioned** resources, managed and scheduled separately
- facility and computing equipment **designed separately**

→ Program = executable on single machine

Warehouse-scale Computer

- computing system designed to run massive internet services
 - e.g., Facebook, Amazon, ...
- applications: tens of binaries running on hundreds to thousands of machines
- **homogeneous hardware and system software**
- **common pool** of resources that are managed centrally
- **integrated design of facility and computing machinery**

→ Program = an internet service

Addressing Network Bottleneck

WSC는 많은 server로 이루어짐

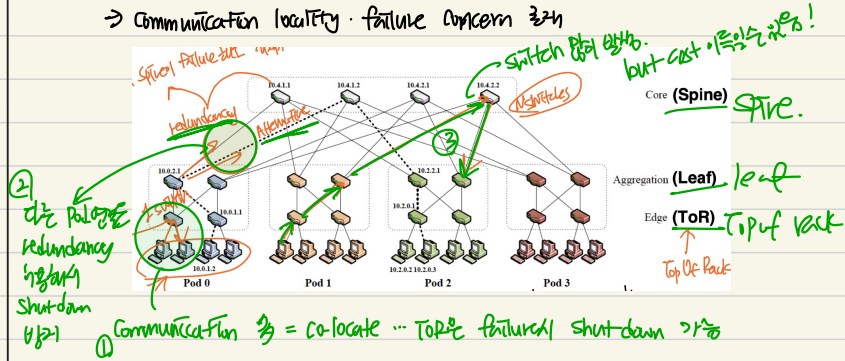
leaf hierarchical network

Bandwidth : Cost per gigabyte port

- 48-port (rack) switch : \$400 (이것을 많이 이용하면?)
- 10,000 port (data center) switch : Priceless

approach for lower-cost

- large-scale networking을 cheap cost로 하기?) 성능은 우선시
- collection of small commodity switches
- communication locality · failure concern 지니



→ network bottleneck & redundancy high

Flash (SSDs) in Warehouse-Scale Machines

Disk

Flash value proposition

Storage hierarchy

- DRAM을 대체하기엔 아직 속도가 없지만, 많이 양산중
- Still slower than data 특성에 따라서 다르게 접근해야함.
 - Cold / warm / hot 데이터 type 나누기
 - Gold data가 무조건 communicate 과다함!
- Complexity of storage hierarchy 는 점점 더 커진다!

Performance Variability · Faults

↓
 Performance Variability
 Scale up/down
 Scale up/down

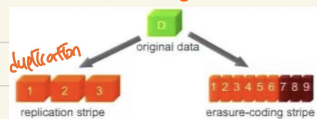
Performance Variability

- large scale에서는 Performance Variability가 크다
- Rare slowdown events can hurt
 - ex) Avg server response time 100ms 일지, every 10k query마다 1s 정도 random performance hiccup
 200ms → 1000ms 정도로 slow down 됨
 → 이 경우, 1 machine = 0.01% slow queries
 5k machine = 50% slow queries ← 최악(?)
- Minimizing performance idiosyncrasies 가 중요하 = rate 할 것 없지 않 (독이네)

Service Level Outages

- Causes
 - networking issues
 - Power Issues
 - ops Issues (automated / manmade mistakes)

- Solution
 - Storage Availability



- duplication 여러개 두는 것
- different disk / server / racks
- distribute data blocks across whole cluster
- load-balanced read 가능
- recovery from faults easy

Demands of GPU clusters

Demands

Resource handling / management

- DL training on GPUs
- high-speed network

Failure handling

Storage - data handling

- sharing identical data > locality 증가하기 중요!
- reusing checkpointed models

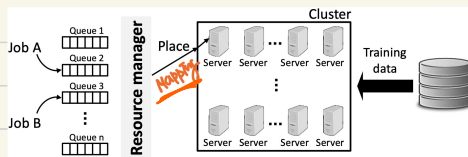
GPU cluster Architecture

- cluster : 100s of servers and thousands of GPUs
- Distributed filesystem : shared storage
- Resource Manager



collective 필요 fairness / efficiency 고려 필요하다.

- Job queue : resource (GPU) amount allowed for each group

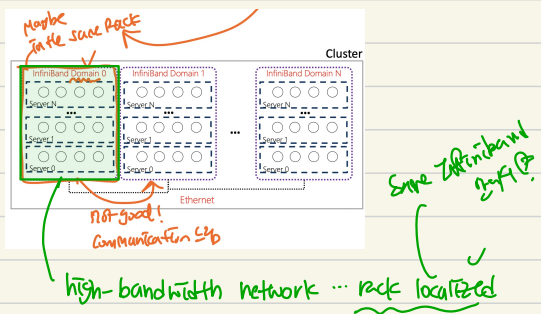


→ idle resource management, load-balancing handling

Demands of GPU clusters (Cont.)

Communication Heterogeneity

- Communication frequency \uparrow \Rightarrow 애플리케이션이 rack-localized 될 수 있도록!

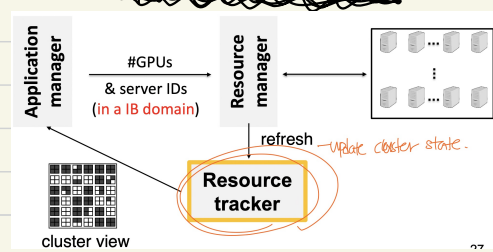


Job Placement

- Job's GPU allocation도 최대한 같은 서버 상에서! \Rightarrow localized 할수록 \uparrow
- " within single InfiniBand domain

Resource Negotiation

- each job은 ~~각각~~ resource requirements에 따라 application manager가 결정.



- Intra-server (MLLB) locality에 따라...

