

How can we crunch the massive amount of data?  $\Rightarrow$  cloud data center

(store and process big data)

$\rightarrow$  94% global workloads was processed in datacenter (computations)

elastic computing cloud (EC2).

Read papers.

The first pass (5 columns): title, abstract, introduce, conclusion.

category, context (which other paper is it related to), correctness, contribution (is it significant).

The second pass (1 hour): content, no details

x proof and implementation  $\checkmark$  figures, diagrams. ? convincing  $\rightarrow$  conclusion

The third pass (several hours and days) technical details.

interested in and have time.

Commodity server  $\xrightarrow{\text{you 80\%}}$  rack  $\longrightarrow$  cell (油压供油).

a GPU server: memory

psb  $\rightarrow$  cpu  
bottle  
neck

GPU  $\rightarrow$  NVLink RDMA used to connect to remote GPUs bypass CPUs

inter-rack: optical circuits switch. based programmable network

<sup>叶级</sup> spine-leaf arch: spine aggregation switch. + distribution switch (fat tree). 每层带宽不同, 靠近根部带宽增加.

GPUs have dedicated interfaces with rack switches which can be used to connect with remote GPUs by pass CPU

users → network → border router → cluster router → layer 2 switch

→ rack switch

Inter-DC WAN is very slow.

普遍存在

ubiquitous, convenient, on-demand network access.

on-demand computing delivered over the Internet with pay as you go pricing. infrastructure as software.

## Lecture 2.

For most of users the best way to deploy large scale internet services is just to deploy on cloud.

VS

Economies of scale. / statistical multiplexing / profitable business.

## Cloud pricing

Computing (EC2), storage (S3, EBS; tiered pricing). Data transfer (inbound free, <sup>数据模式的。</sup> outbound x <sup>data gravity price.</sup>)

1. Reserved pricing:  $cost(t) = U + discount \times R \times t$

Guaranteed availability  $\downarrow$  一次性的费用  
按需价格的常规费用

2. Spot pricing.

user's bid  $\rightarrow$  spot price, win until a new one is posted.  $\rightarrow$  running users with a lower bid get their instances terminated.  
e.g. machine learning

on demand service guarantee: if instance is launched, it will not be terminated by platform  
It is possible the price of (on demand) is lower than (spot), because cloud provide want to

pick out all spot users.

\* brokerage service. 中间服务: make instance acquisition strategies.

Private cloud can be cheaper than public cloud e.g. GPU cloud.  
public cloud  
hybrid cloud.

Cloud computing stack  
application.  
platform  
infrastructure e.g. AWS  
virtualization  
hardware

Storage  
network

cloud service model

flex Infrastructure as a Service: VMs, computation, storage, network I.  
platform : SDK, API, automatic scalability; x control of OS etc.  
software : office365

base Other: Function: write apps in the form of cloud functions and define events I only when  
tools. machine learning: pipelines tools functions running I  
model : model api (chatgpt).

e.g. AWS lambda.

Issues of cloud.

Availability  
data loss

Vendor lock in : lock into the current provider.

↳  
security  
privacy.

challenge: storage, scale. faults

failures. networking