



Use Tesla to provide first GPU VM Service in China

Feng Zhu

专注·服务·中立

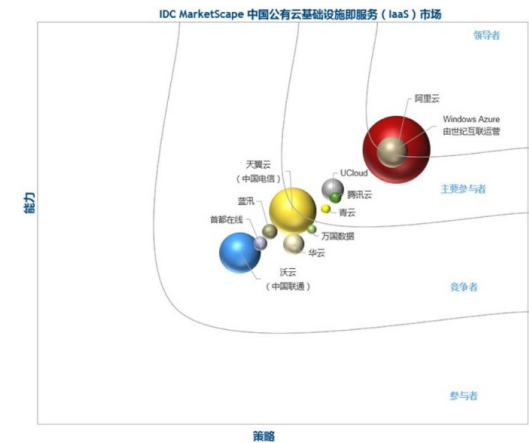
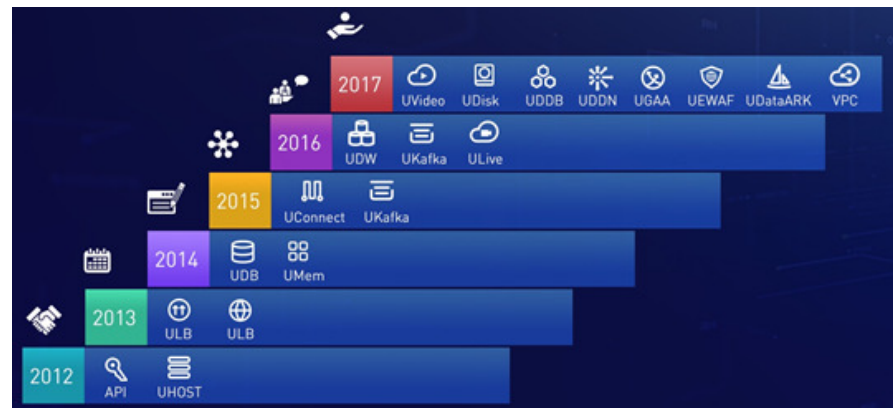


Outline

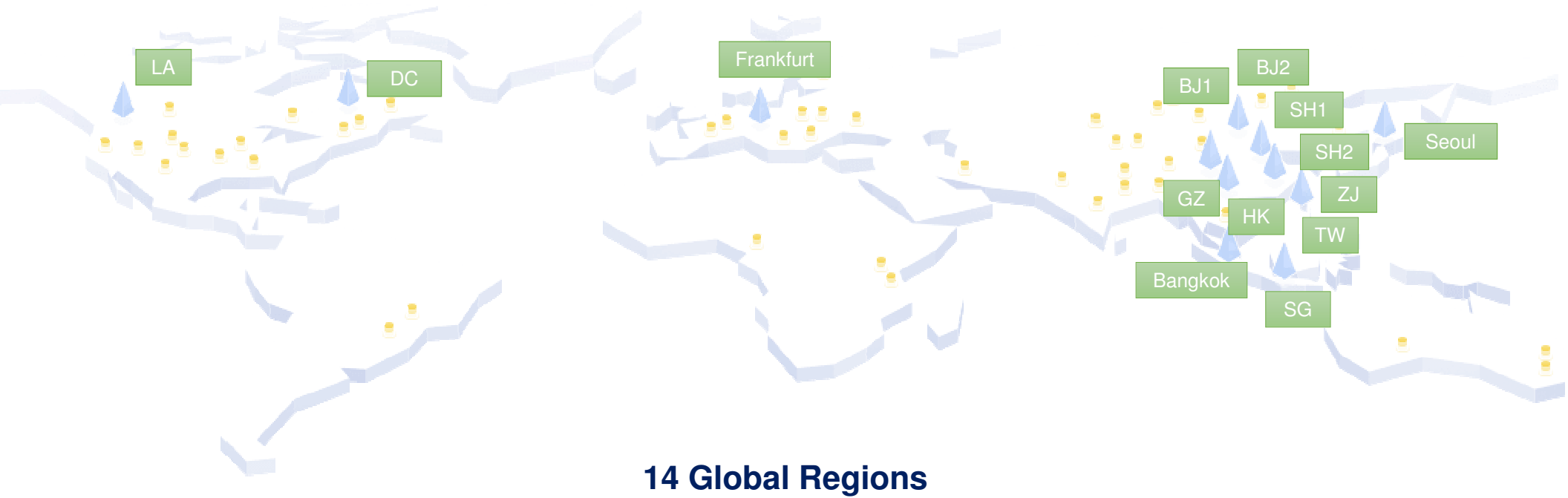
- **UCloud Introduction**
- **K80 GPU VM**
- **P40 GPU VM**
- **UCloud GPU PaaS Service: UAI-Service**
- **UCloud GPU ecosystem**

About UCloud

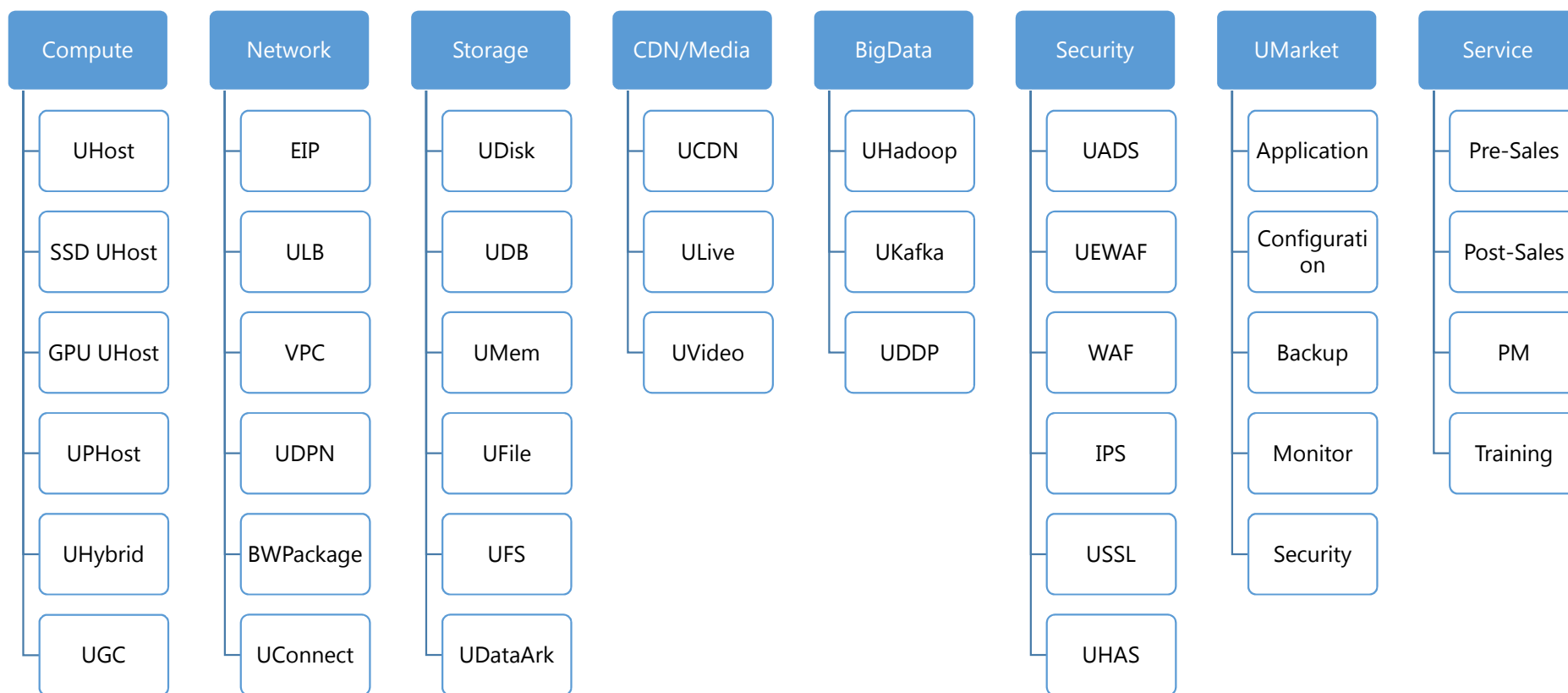
- Top 3 IaaS Provider in China
- Found in 2012
- HQ in Shanghai
- Served 50,000+ Enterprise



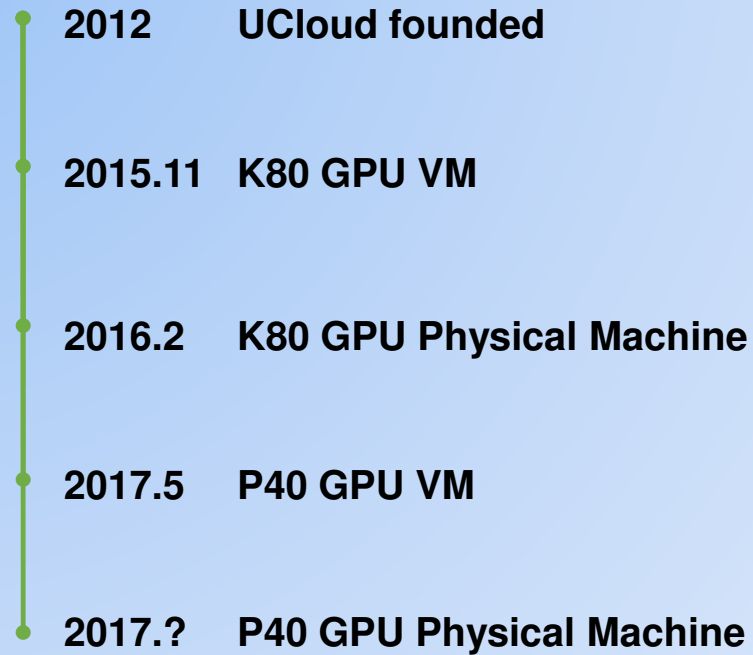
Data Centers



UCloud Product Line



GPU Timeline



GPU Decision: Virtualization

PCI Pass-through	Grid
Single VM occupy GPU device	GPU device shared by multiple VMs
Performance guaranteed	No guarantee on performance
Limited by PCI device number (K80 = 2 * GPU)	no limit for PCI device
No license	License
√	X

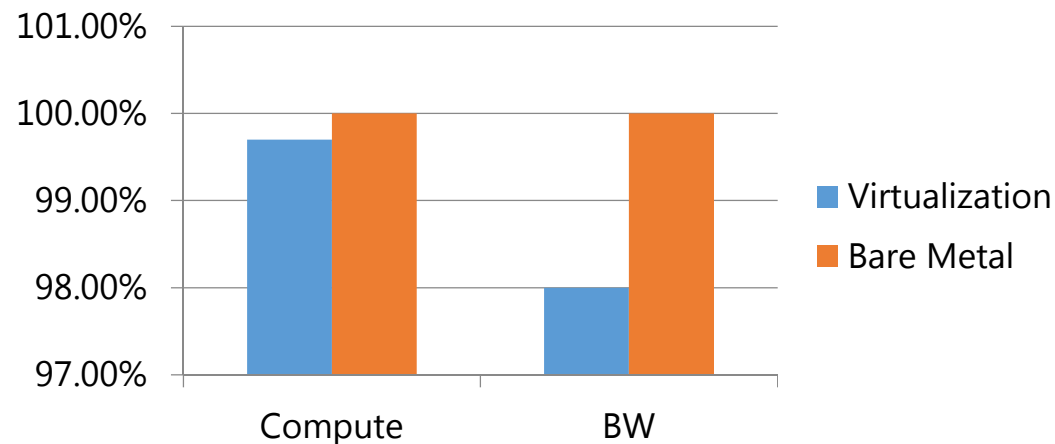
VM Advantage

- **Flexibility for VM configuration**
 - CPU、Memory、Disk size、GPU number are all flexible
 - SDN network flexible
- **Main OS all supported, Win/Linux**
 - CentOS 6.5/CentOS 7.0/Ubuntu 14.04/Ubuntu 12.04/Gentoo 2.2/Win 2008/Win 2012
- **Fast Deployment**
 - Based on self-defined image, can deploy 1000 VMs in 1 minute

VM Performance Degradation

- **Using Pass-through Technology, almost no degradation**

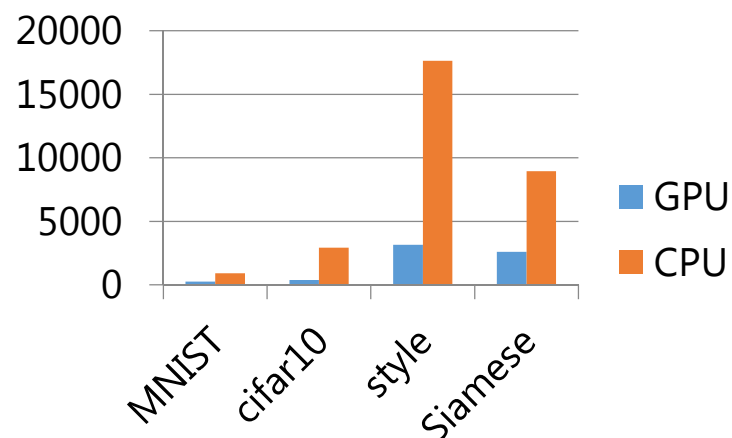
Degradation	Virtualization	Bare Metal
Computing	99.3%	100%
Bandwidth	98%	100%



UCloud GPU Virtualization – DL test

- **Caffe Performance (Ubuntu)**

Cases	iters	GPU(secs)	CPU(secs)	Speedup
LeNet on MNIST	10000	253.9	900.3	3.5
cifar10	5000	374.2	2931.0	7.8
Fine-tuning for style recognition	10000	3138.9	17634.8	5.6
Siamese Network Training	50000	2584.2	8937.1	3.5



UCloud GPU Virtualization – DL test(2)

- **Theano/Keras (Ubuntu)**

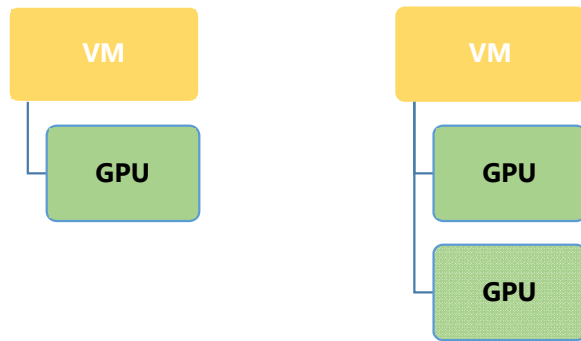
Cases	iters	GPU(secs)	CPU(secs)	Speedup
imdb_lstm.py	20000	47	239	5.1
imdb_cnn.py	20000	11	563	51.2
imdb_cnn_lstm.py	20000	27	236	8.7
addition_rnn.py	45000	5	32	6.4
cifar10_cnn.py	50000	198	2670	13.5
babi_rnn.py	950	9	22	2.4
mnist_cnn.py	60000	23	1332	57.9
mnist_irnn.py	60000	270	451	1.7
mnist_mlp.py	60000	1	5	5.0

K80 Physical Machine

Hardware	Specification
GPU	Tesla K80
CPU	Intel E5 2630*2
Memory	192G
Disk	2T SSD
Networking	10G NIC*4



VM Configuration - K80



CPU

4C

8C

16C

Memory

8G

16G

32G

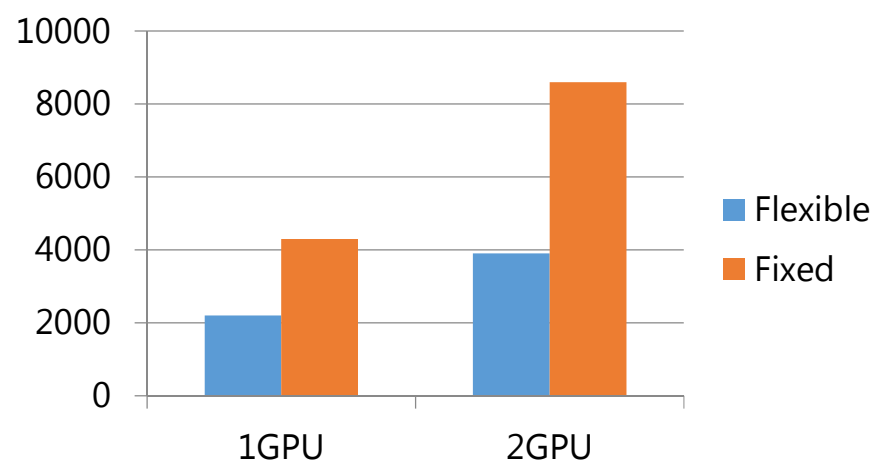
64G

Disk

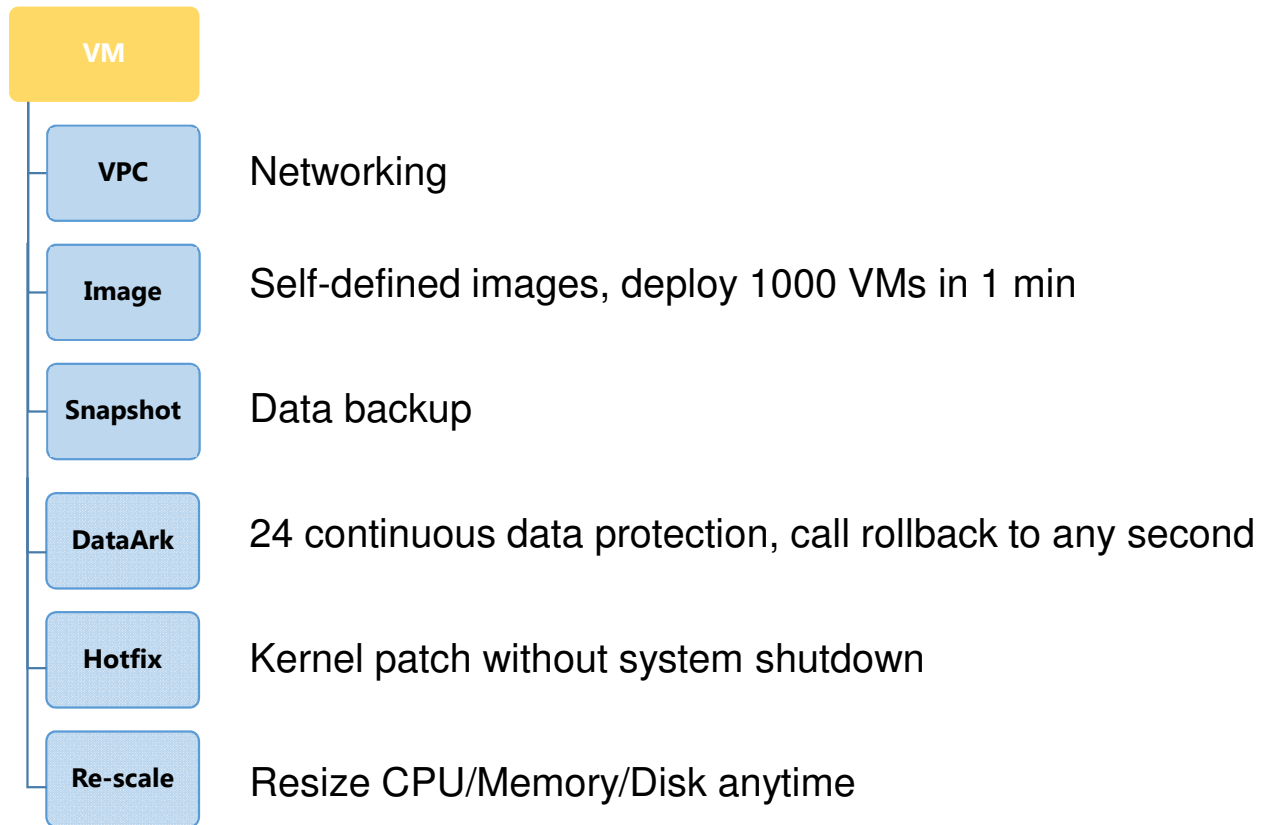
100G - 1T

Flexible VM Save Cost

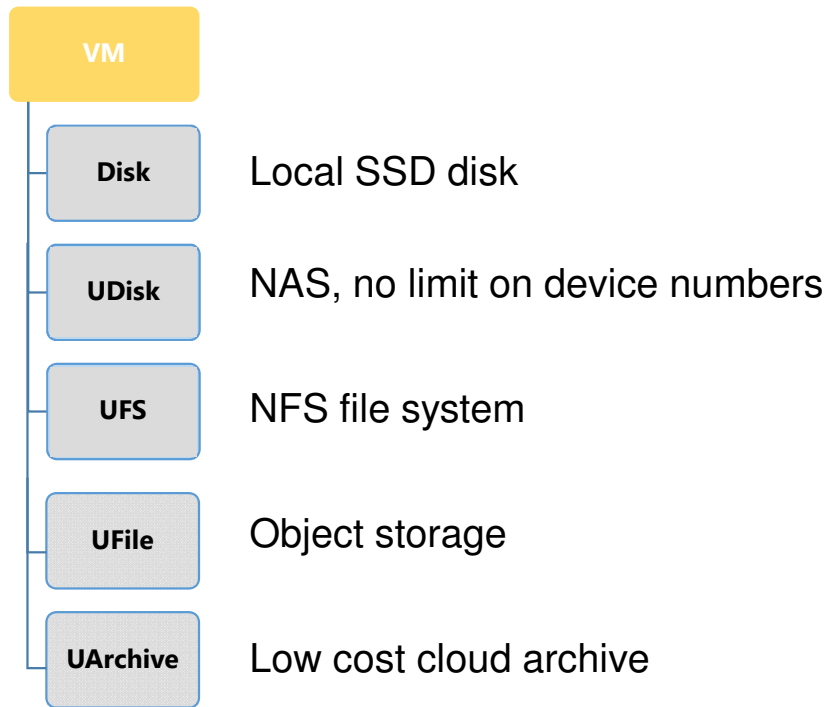
Configuration	Fixed	Flexible
CPU	16C	4C
Memory	96G	8G
Disk	1T	100G
GPU	1	1
Price	4,300 CNY/month	2,200 CNY/month
USD Price	\$615/m	\$315/m



GPU VM Features



Storage Solution



Create GPU VM

U 产品与服务 PD 北京二 全部 uhosttest@ucloud.cn

云主机 UHost 切换旧版

主机管理 镜像管理 回收站

创建主机 启动 关闭 更多

主机名称	资源ID	业务组	基础网络	配置	机型 特性	创建时间	状态	操作
云主机 (UHost) 提供可随时扩展的计算服务, 每台云主机以虚拟机的形式运行, 主机资源包含CPU、内存、磁盘等最基础的计算组件。创建主机								

Create GPU VM

U 产品与服务 PD

uhosttest@ucloud.cn

创建主机

地域

北京一	北京二	浙江	上海一	上海二	广州
香港	洛杉矶	华盛顿	法兰克福	曼谷	首尔
新加坡	台湾				

可用区

可用区B	可用区C	可用区D
------	------	-------------

基础配置

系列 **系列1**

机型 **GPU型G1** [查看介绍](#)

网络增强 OFF 获得数倍于普通机型的网络包处理性能

购买数量

云主机配额 1/994

弹性IP配额 1/378

付费方式 月付

购至月末

合计费用 **1,554.95元**

[立即购买](#) [加入清单](#)

Create GPU VM

U 产品与服务 PD

uhosttest@ucloud.cn

创建主机

系列 ① **系列1**

机型 **GPU型G1** 查看介绍

网络增强 ② OFF 获得数倍于普通机型的网络包处理性能

CPU **4核** 8核 16核

GPU **1颗** 2颗

内存 **8G** 16G

镜像 ③ **标准** 自制

CentOS CentOS 6.5 64位

存储类型 ④ **本地SSD盘** 云硬盘

系统盘 20GB 创建后可通过“更改配置”操作扩容

数据盘 **250GB** 500GB 750GB 1000GB 100

购买数量

云主机配额 1/994

弹性IP配额 1/378

付费方式 月付 购至月末

合计费用 1,554.95元

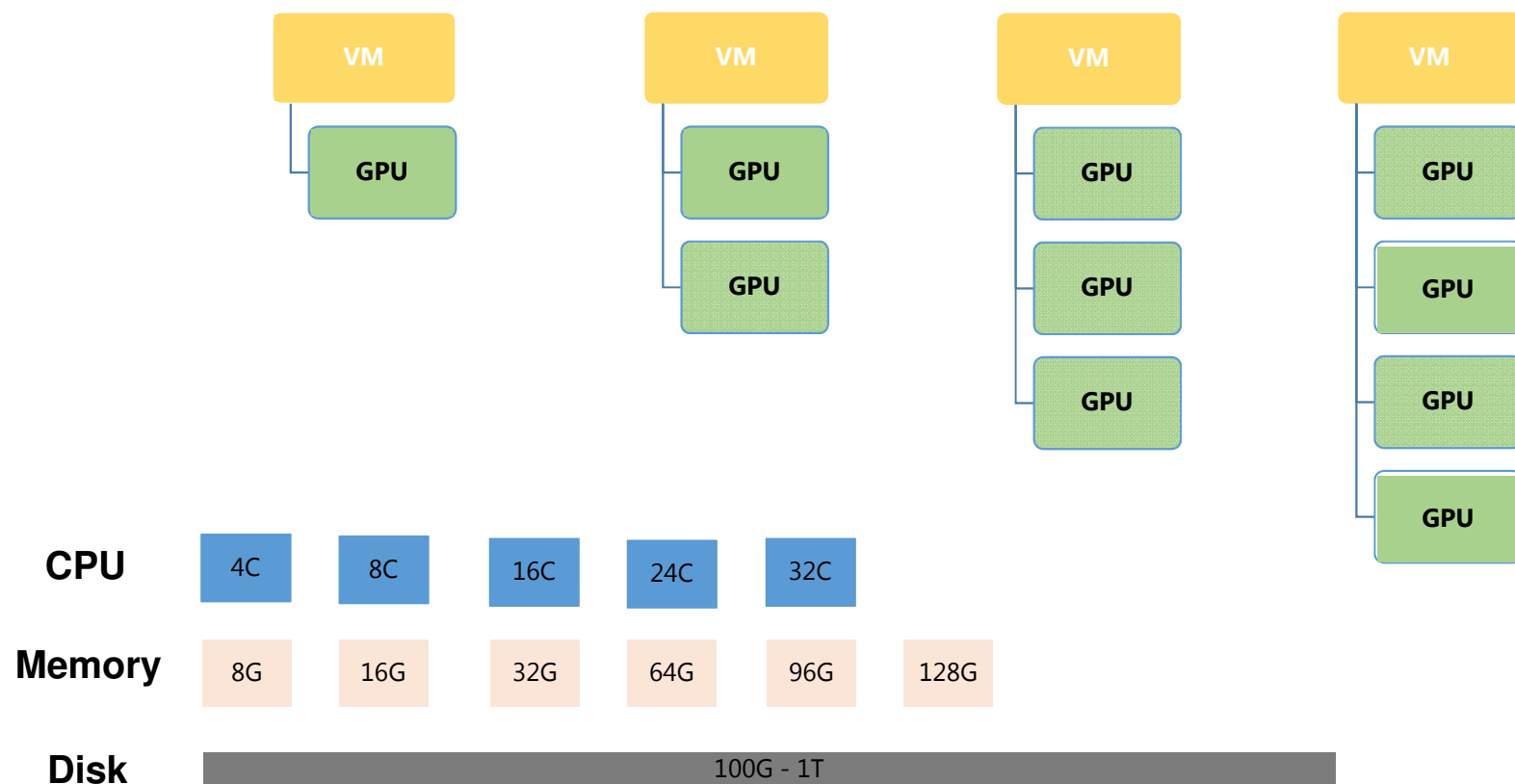
立即购买 加入清单

P40 Physical Machine

Hardware	Specification
GPU	Tesla P40*4
CPU	Intel E5 2650*2
Memory	256G
Disk	3T SSD
Networking	10G NIC*4

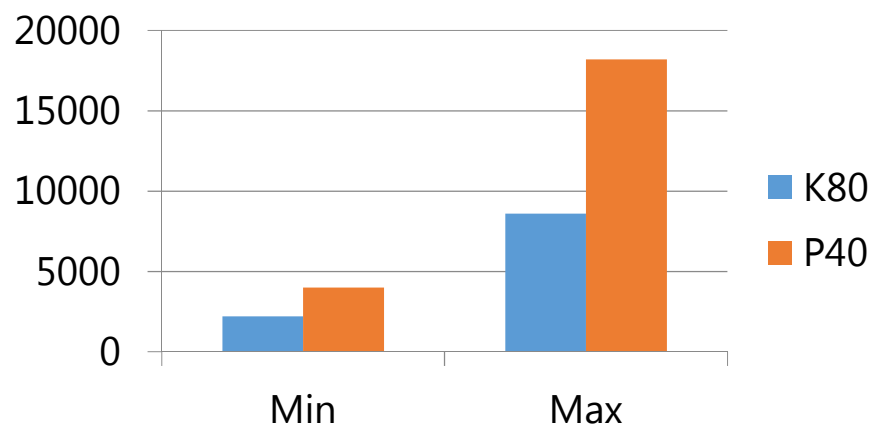


VM Configuration - P40

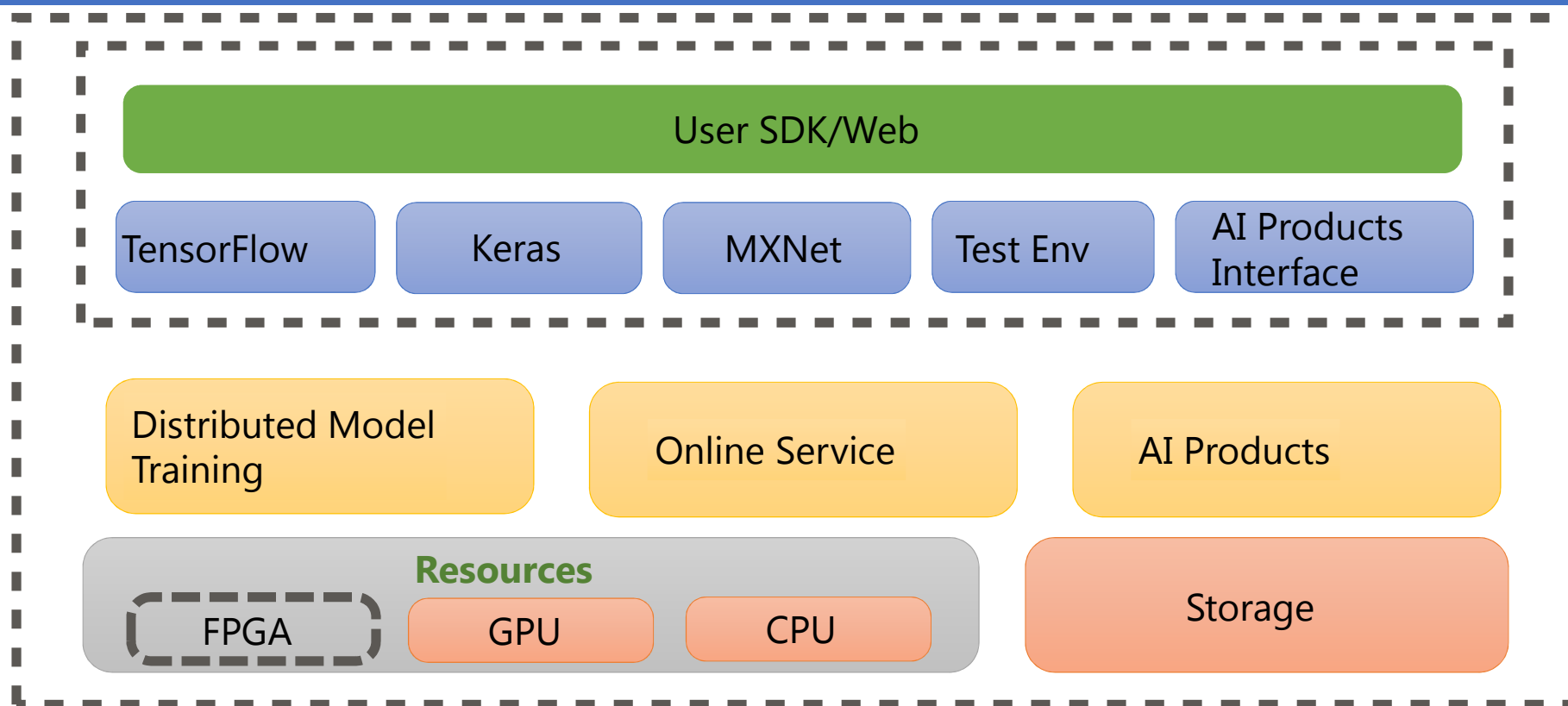


P40 Price

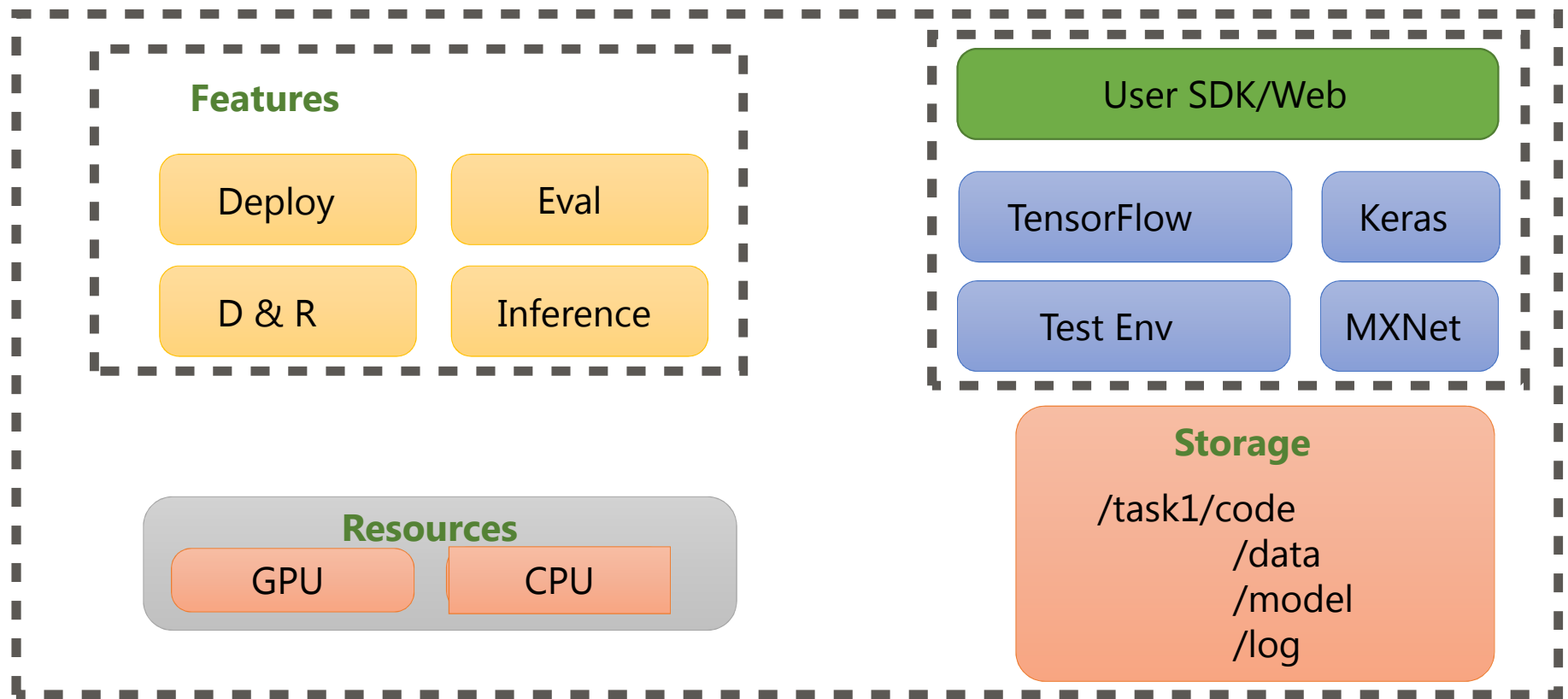
Configuration	Spec 1	Spec 2
CPU	4C	32C
Memory	8G	128G
Disk	100G	1T
GPU	1	4
Price	4,000 CNY/month	18,200 CNY/month
USD Price	\$570/m	\$2,600/m



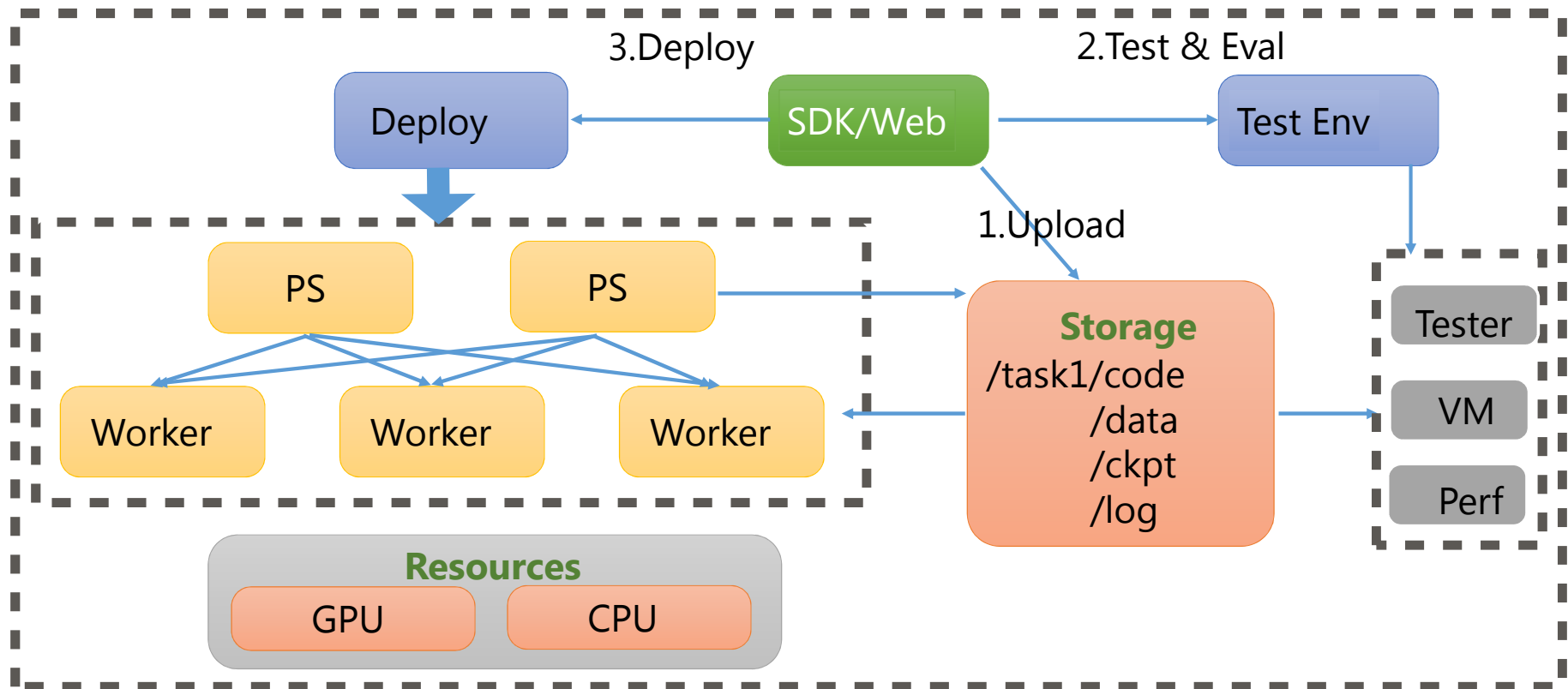
UAI-Service Overview



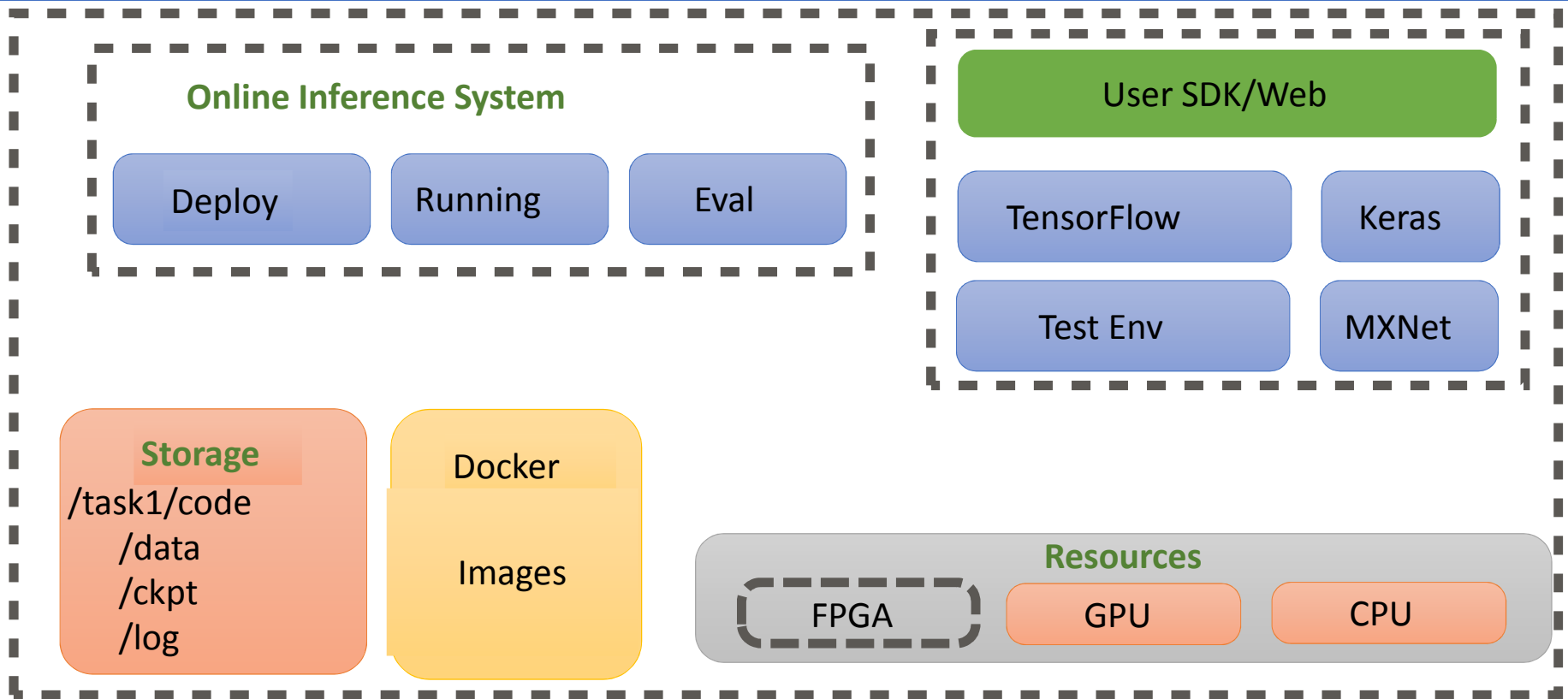
Distributed Training Layout



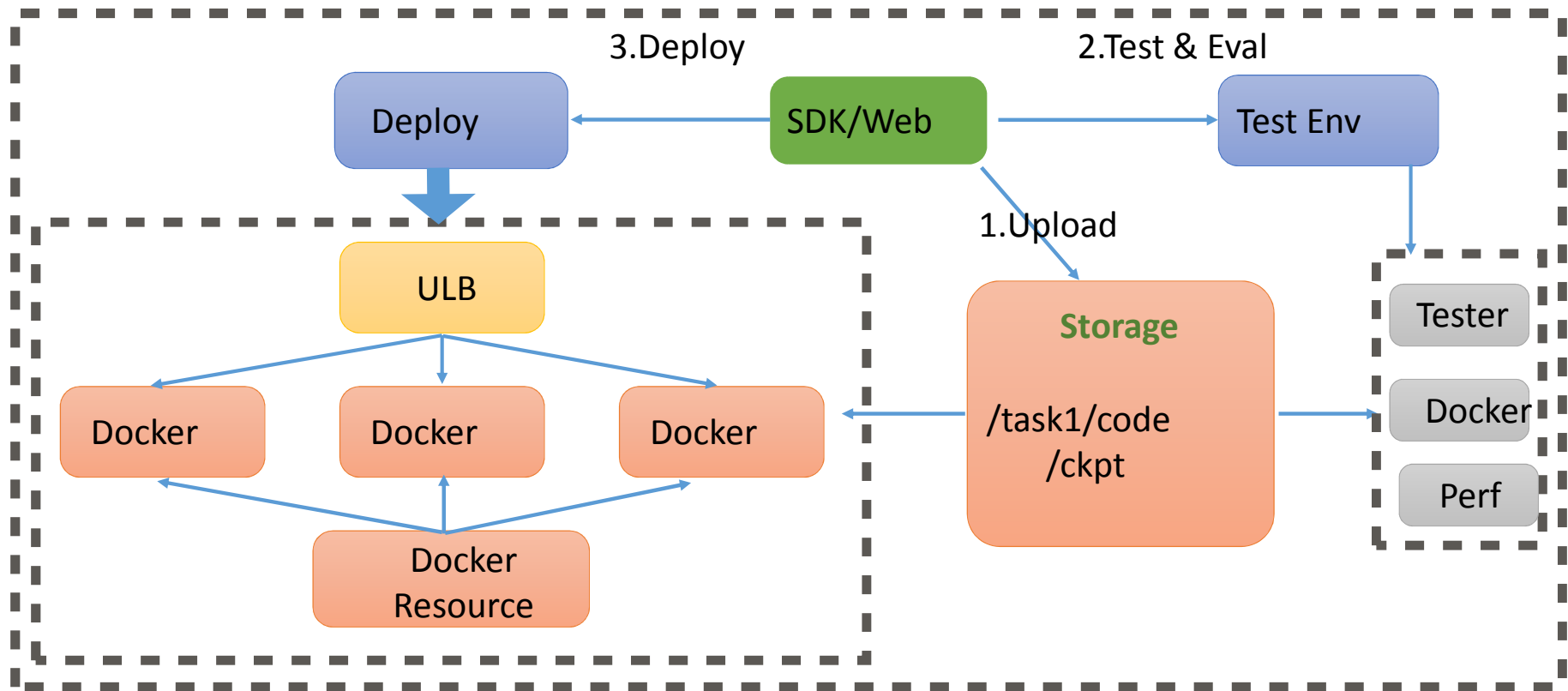
Distributed Training Process



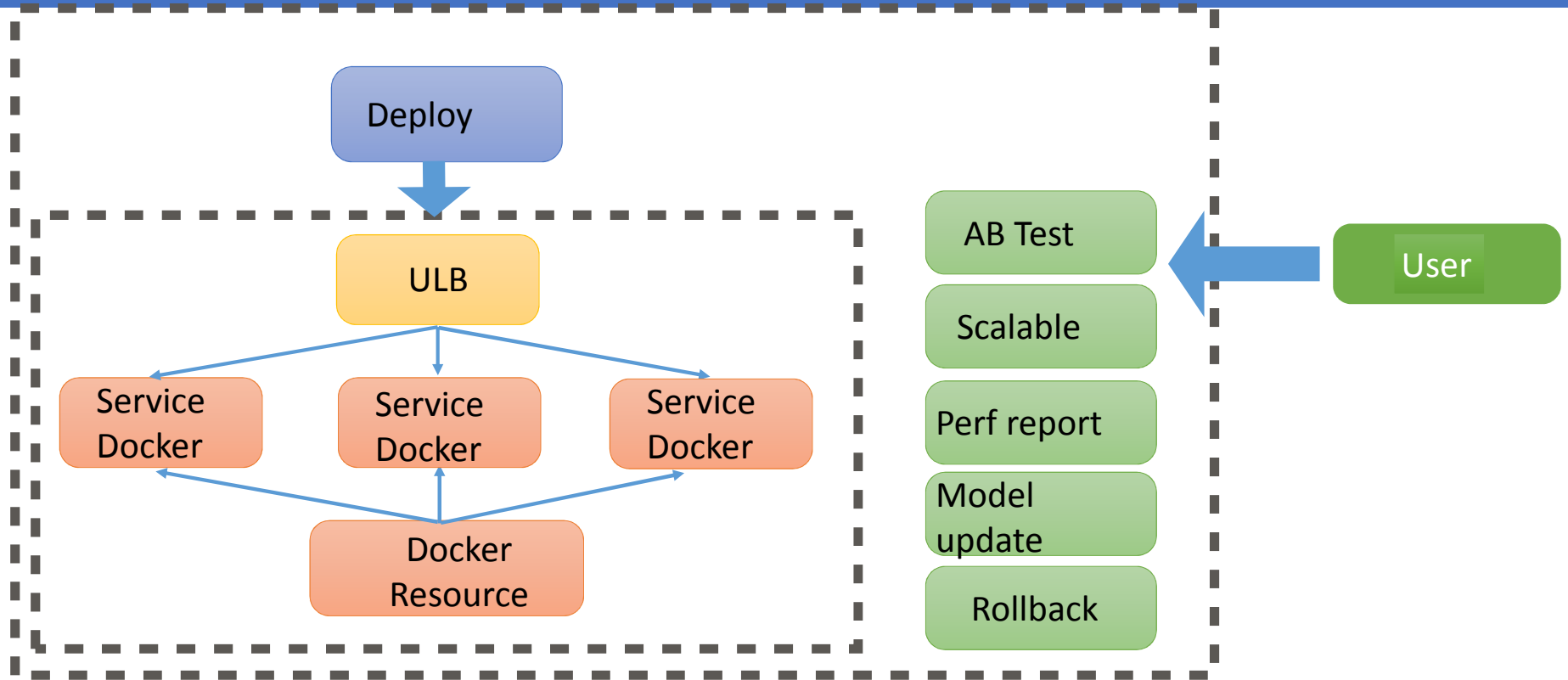
Online Inference Layout



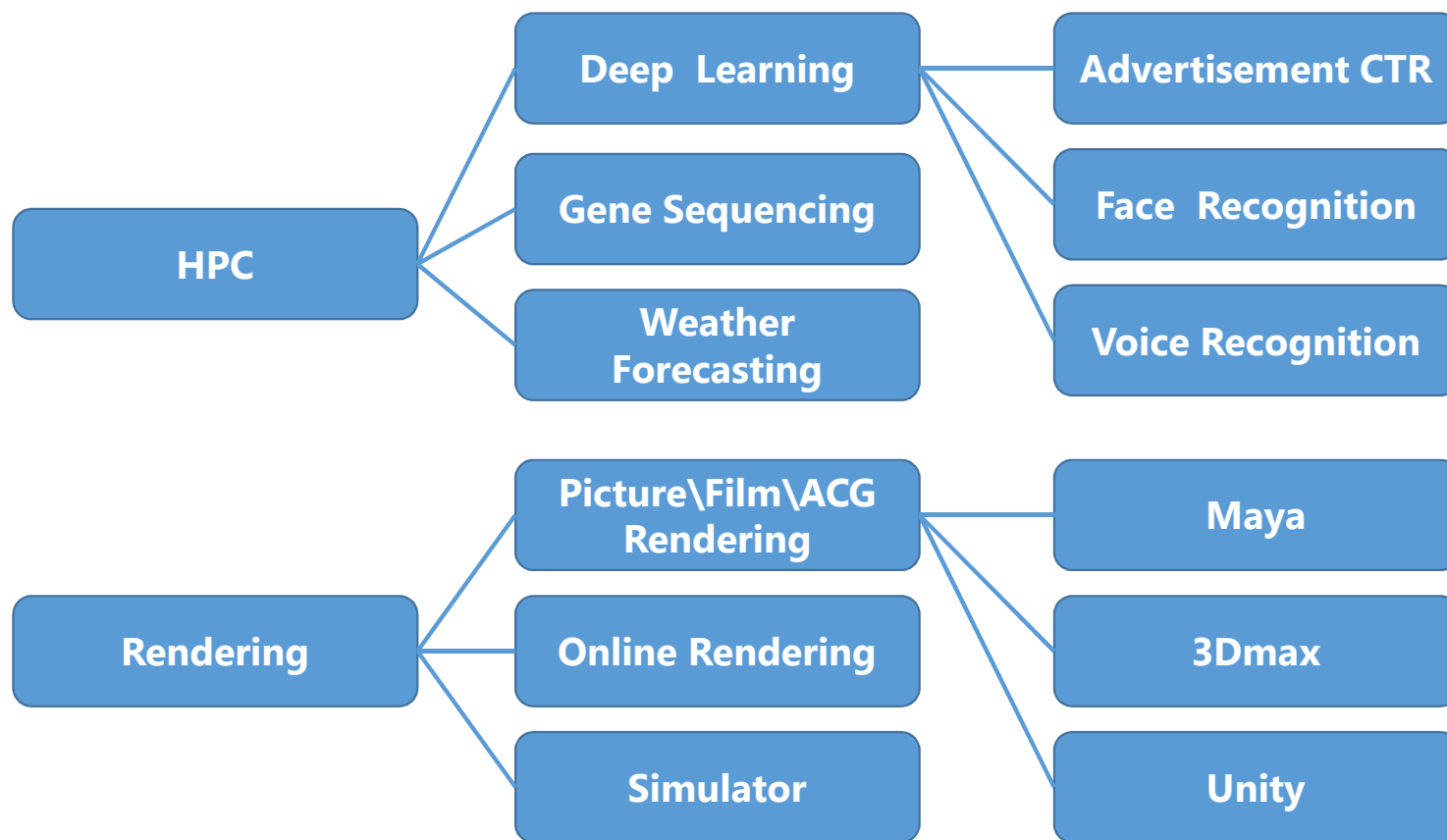
Online Inference Process



Online Inference API/SDK

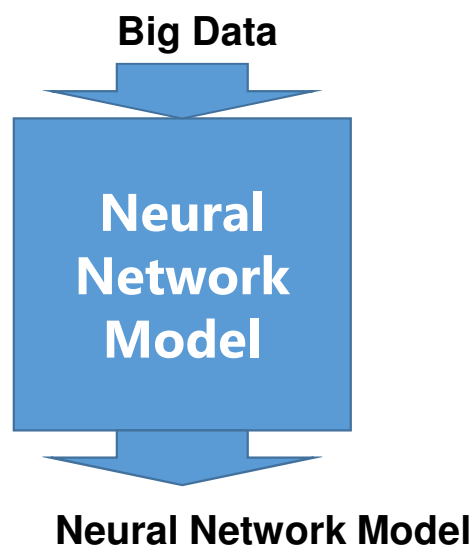


GPU Scenario

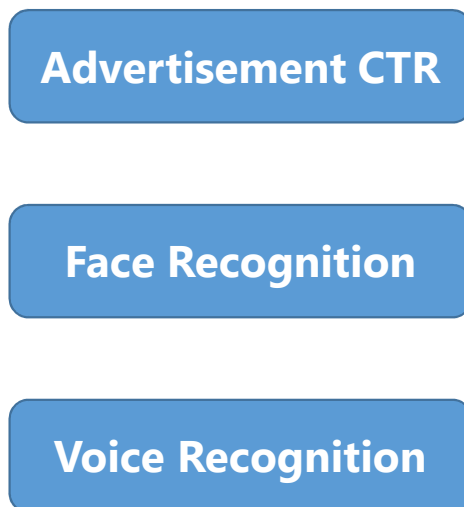


GPU Scenario

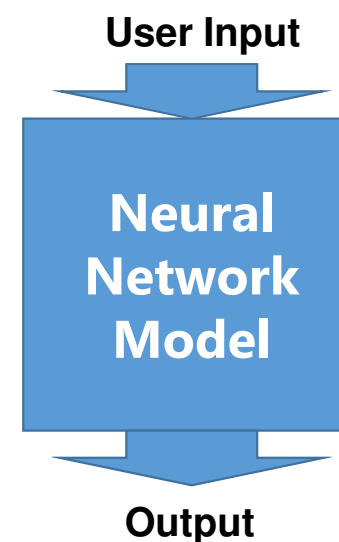
Training



Compute-Intensive
(GPU)



Online Service



Compute-Sensitive
(GPU)

GPU Scenario: Example

- **CTR click-through-rate estimation**
 - Precise Advertise: according to personal historical data
 - Ad bidding & ranking : $\text{CTR click rate} * \text{Bid Price}$
 - Video、UGC recommendation : watching ratio

GPU Scenario: Example

- CTR click-through-rate estimation

$x = [\text{Weekday=Wednesday, Gender=Male, City=Shanghai}]$

$x = [0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, \dots, 0]$

CTR Estimate Model

Percent of Click: 25%

Thank You

www.ucloud.cn

UCLLOUD

[首页](#) [产品](#) [服务](#) [行业](#) [保障](#) [资讯](#) [生态](#) [关于](#)

[控制台](#) [备案](#) uhosttest@ucloud.cn [退出](#)

 Think
in
Cloud
BEIJING国贸三期 2017.03.29

云筑梦想·应需而为

