# Applied AI Architecture @ Alibaba Infrastructure
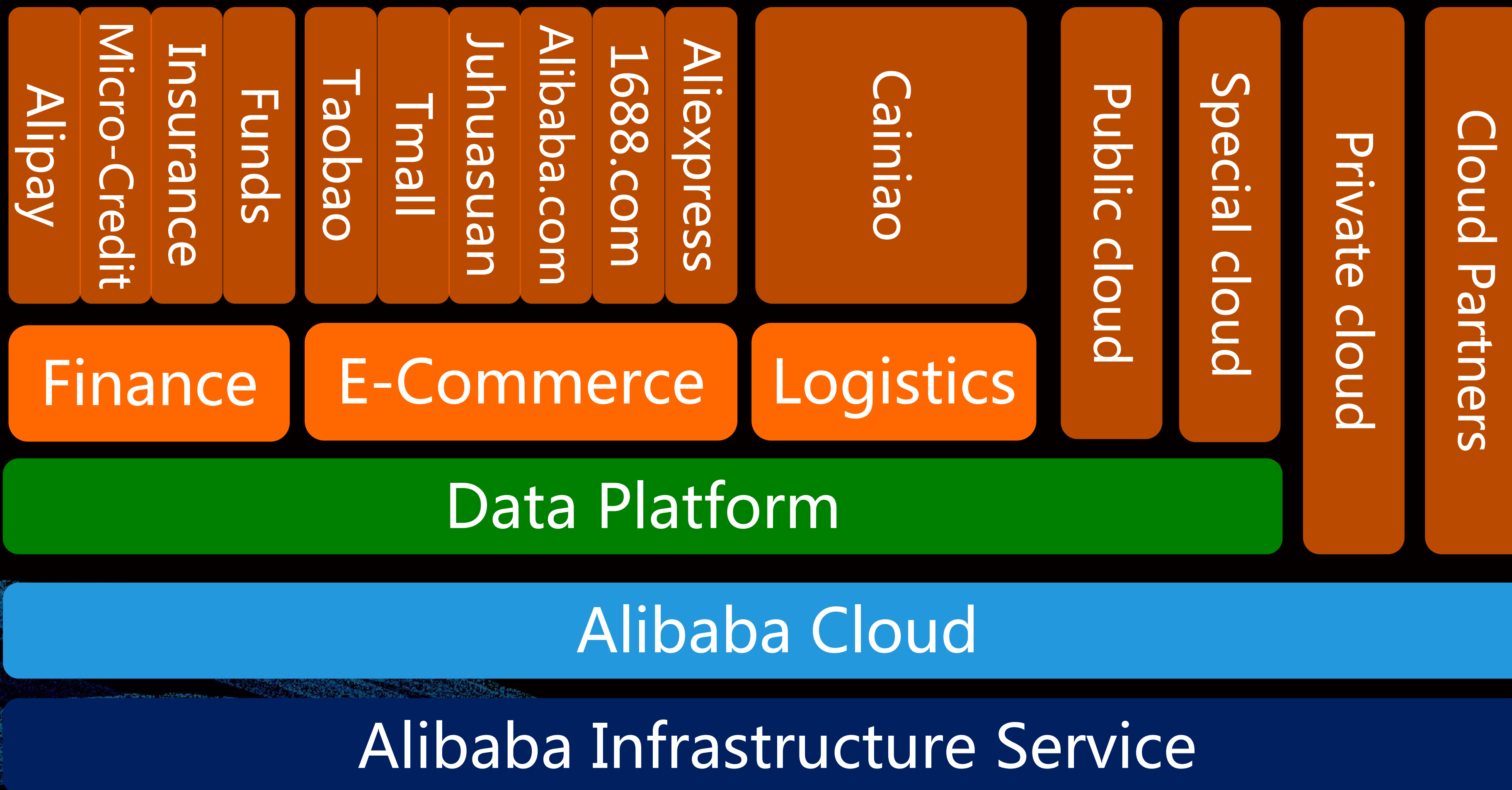
Lingjie XU, Director
Heterogeneous Computing
AIS

# Who I am

- ➤ Joined Alibaba Spring 2017

- ➤ Leading Applied AI Architecture team, focusing on AI HW acceleration and HW/SW co-play

- ➤ Held multiple senior architect and management roles in GPU domain

# Who we are

Alipay | Micro-Credit | Insurance | Funds | Taobao | Tmall | Juhuasuan | Alibaba.com | 1688.com | Aliexpress | Cainiao | Public cloud | Special cloud | Private cloud | Cloud Partners

**Finance** | **E-Commerce** | **Logistics**

**Data Platform**

**Alibaba Cloud**

**Alibaba Infrastructure Service**

# Alibaba Infrastructure



**IDC**

Modular
Eco-Friendly
Automation

**Server**

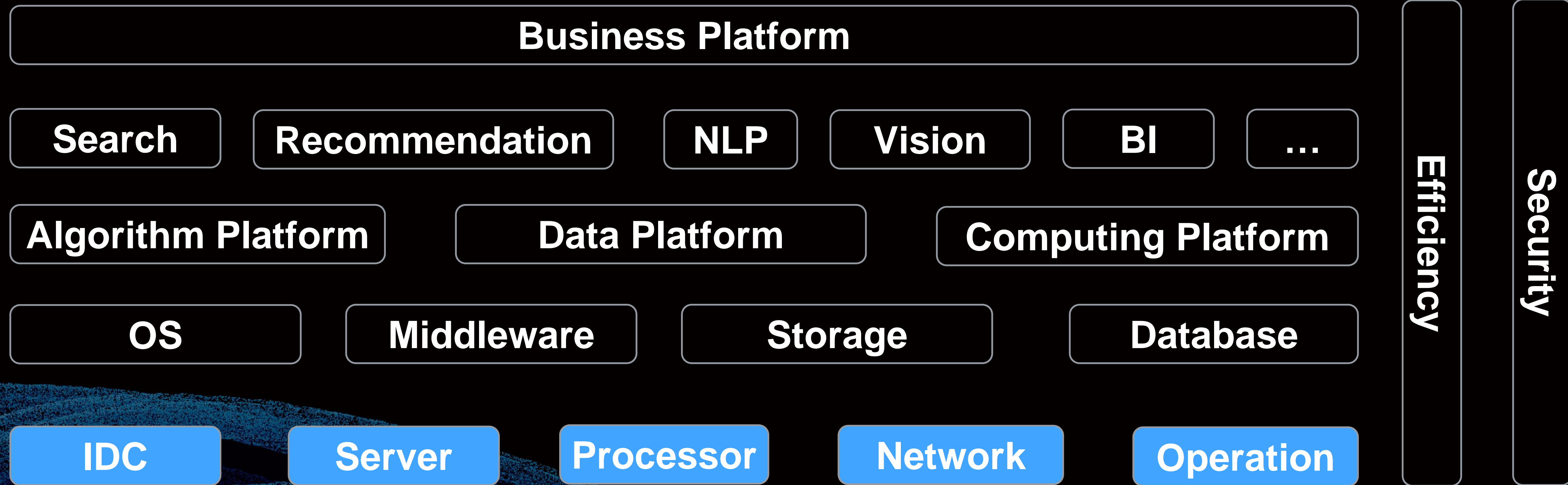High Perf
Low Power
Scalability

**Network**

100G
SDN
Security

**GOC**

Monitor
Analyze
Act

@ 2018 Alibaba Group

# Technology Overview

**Business Platform**

| Search | Recommendation | NLP | Vision | BI | ... |

**Algorithm Platform** | **Data Platform** | **Computing Platform**

**OS** | **Middleware** | **Storage** | **Database**

**IDC** | **Server** | **Processor** | **Network** | **Operation**

**Efficiency**

**Security**

@ 2018 Alibaba Group

# Datacenters



Zhangbei Datacenter
(Fresh air cooling system)
**Best PUE <1.2**
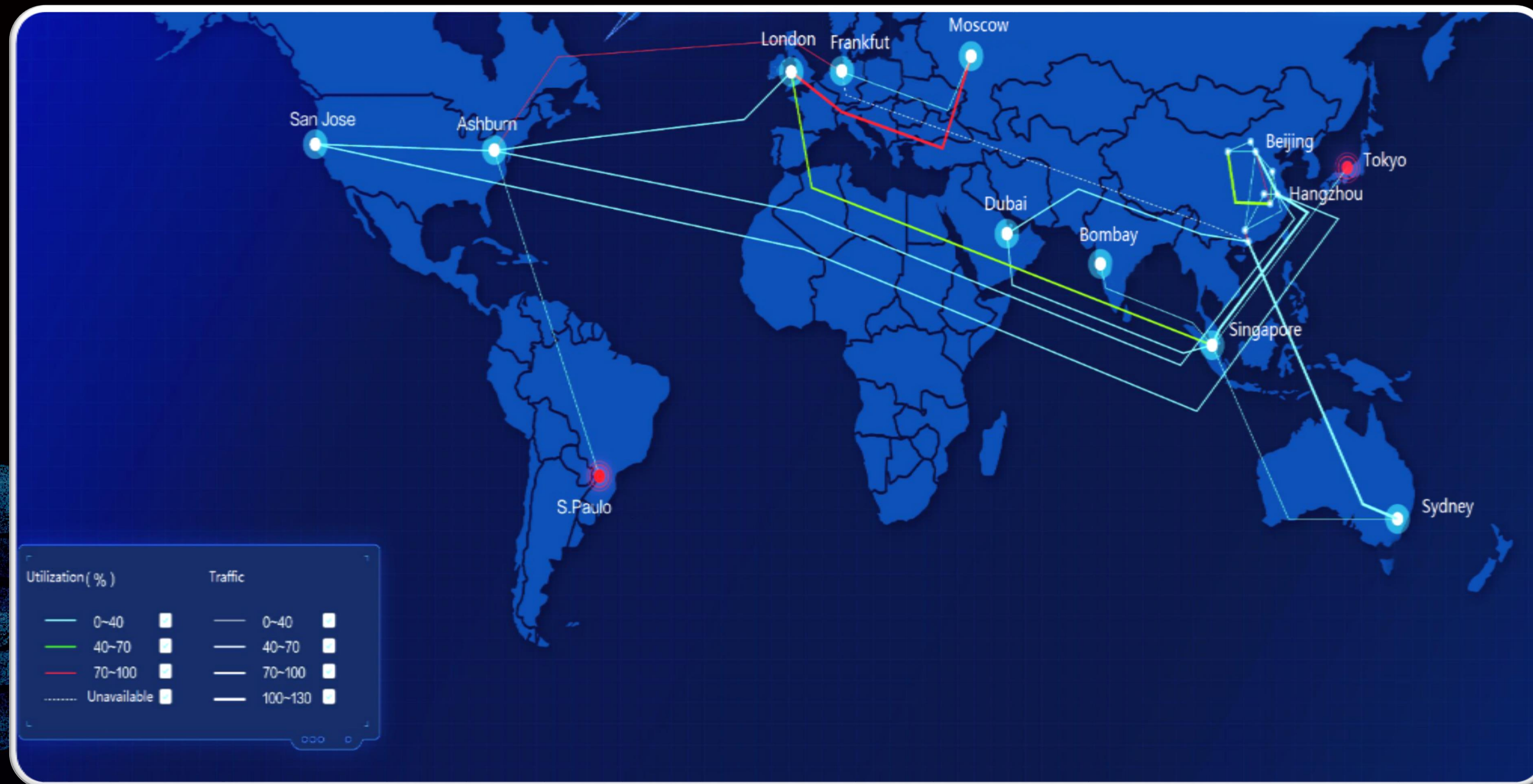
New Frontier
Server immersion cooling
**PUE ~1.0**

Qiandaohu Datacenter
(Lake water cooling system)
**PUE < 1.3**
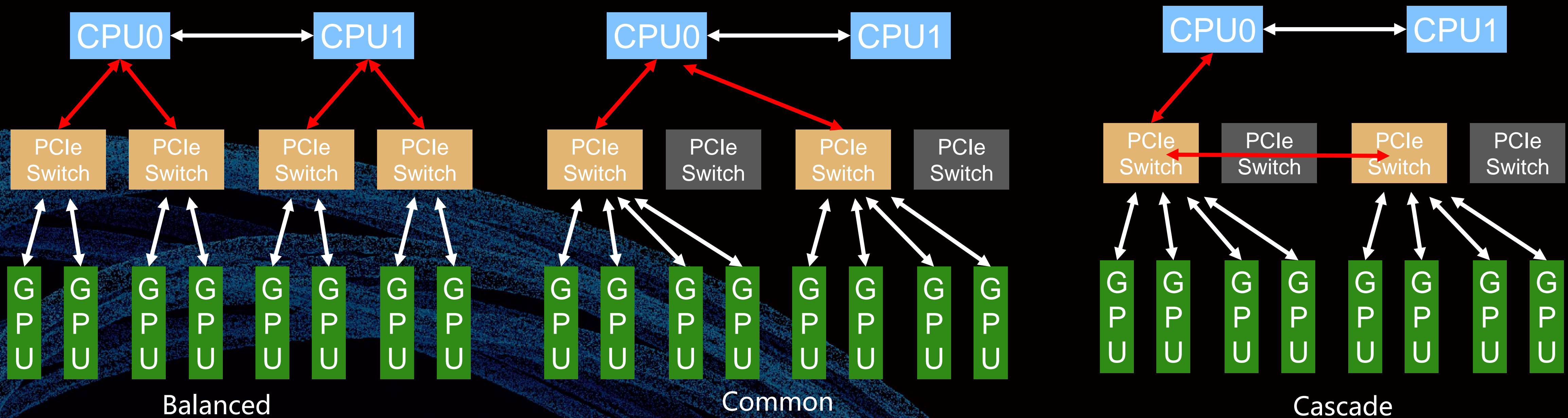
# Network

Massive Scale + Diverse Applications + Bursty Traffic + Fast Growth

# Compute & Storage

# GN6: 8-way GPU Server

- SXM2 or PCIe
- Decoupled modular design
- Configurable topology

Balanced

Common

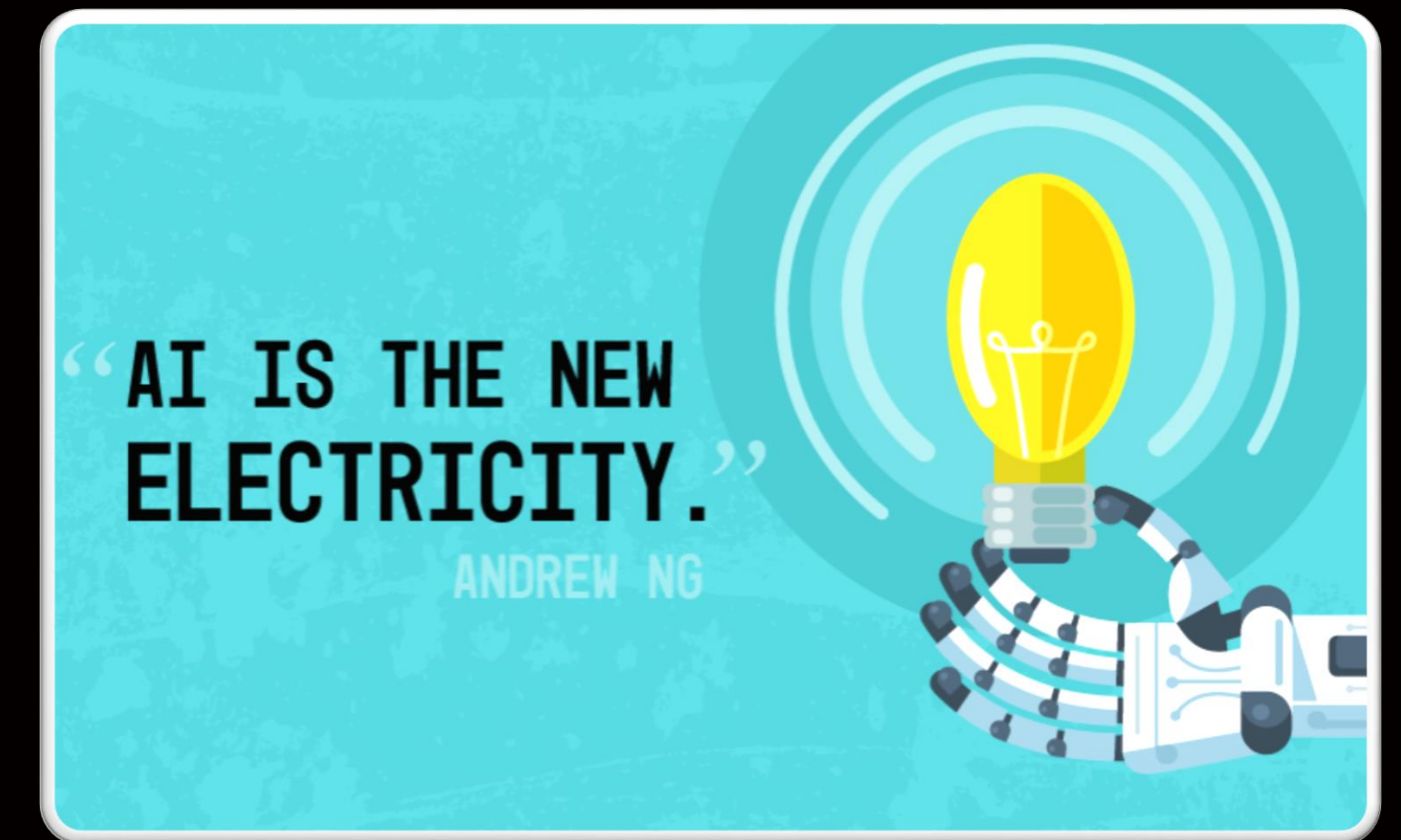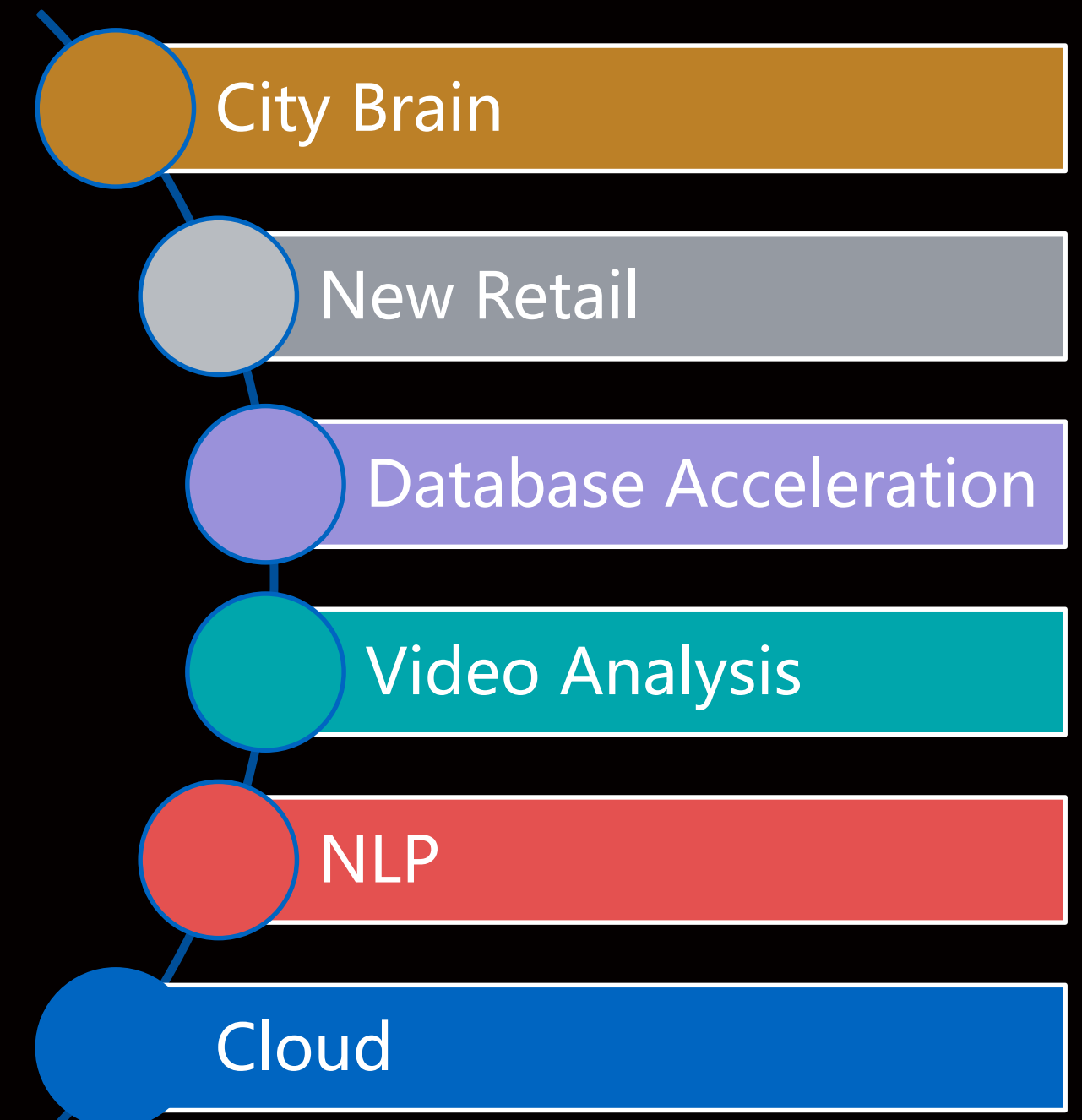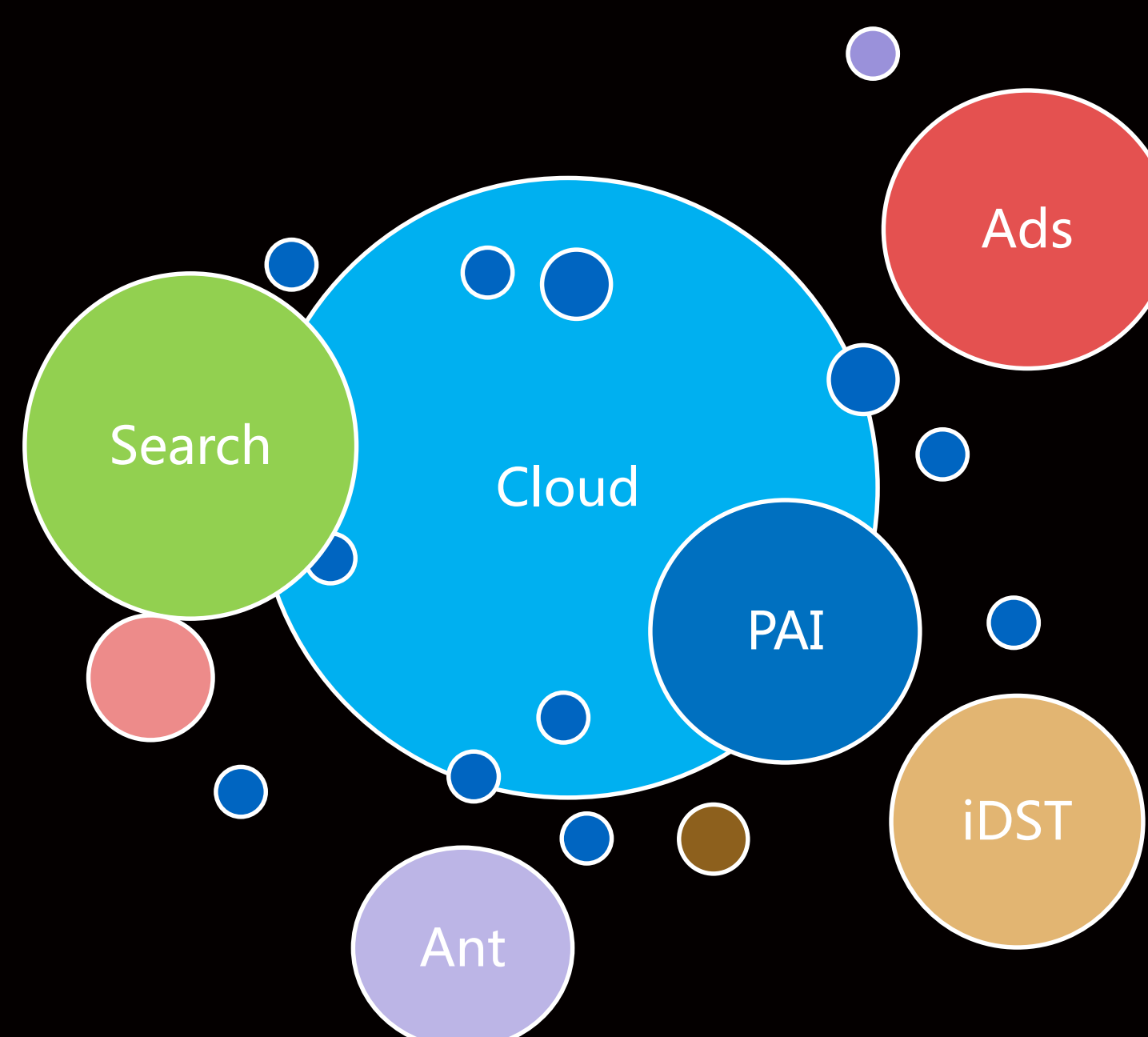Cascade

@ 2018 Alibaba Group

# PaiLiTao

- Category Prediction
- Object Detection
- Feature Extraction
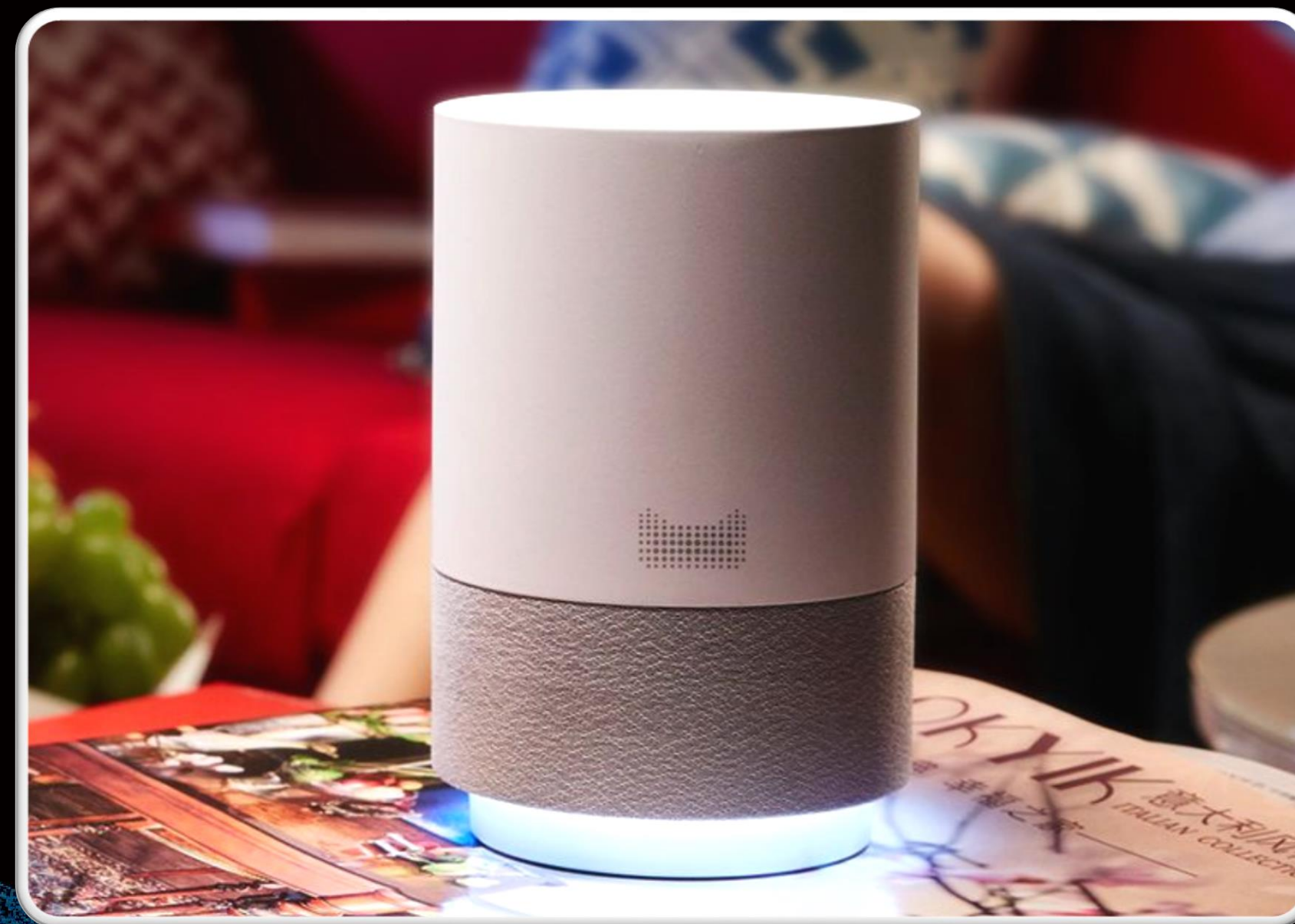- Index Searching
- Soring & Output

# OCR



- 10s Millions of Image
- CNN Model
- Single character accuracy 99.6%
- Overall accuracy 93%
- 8 way distributed GPU solution
- 7x training speed

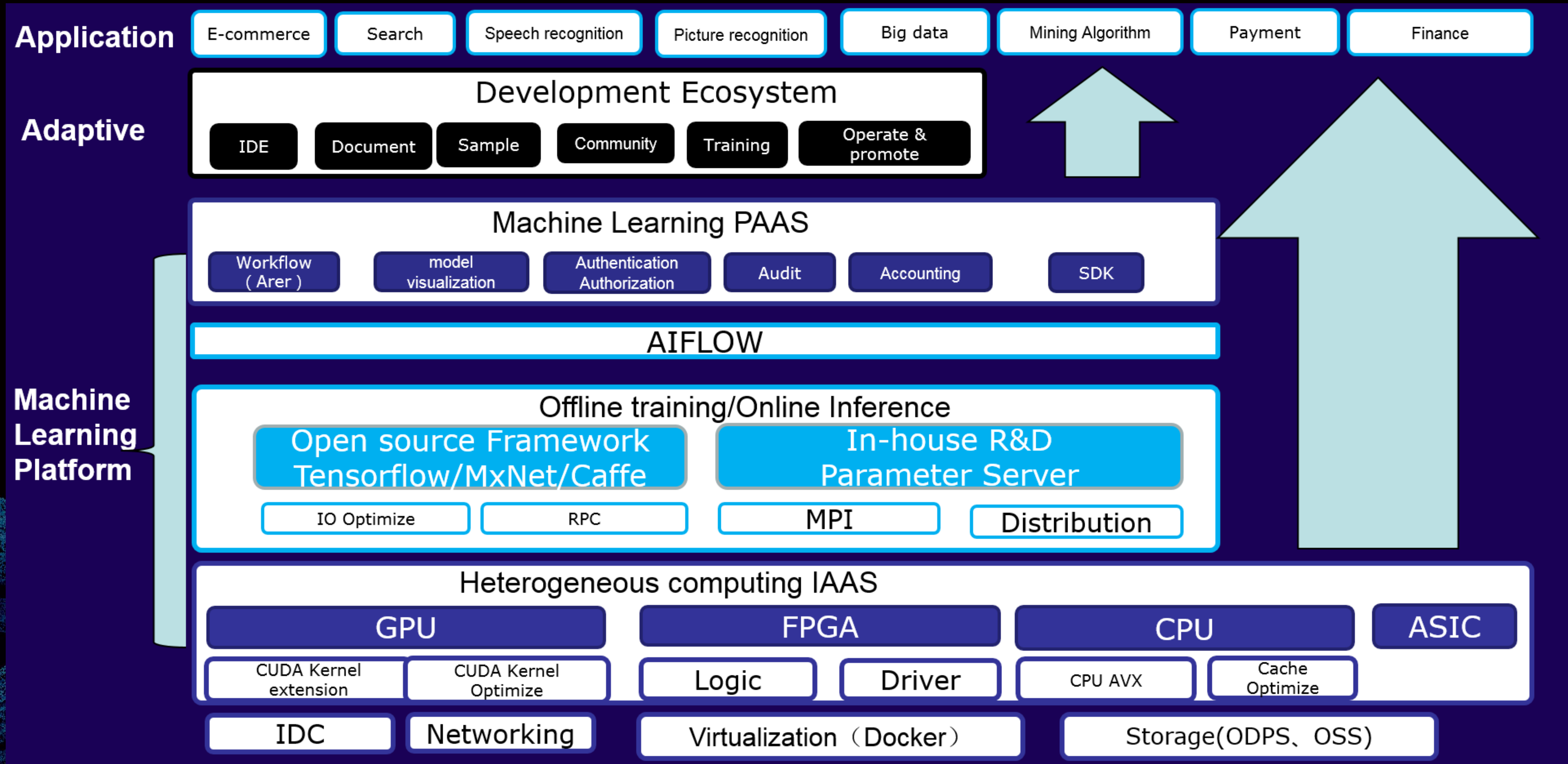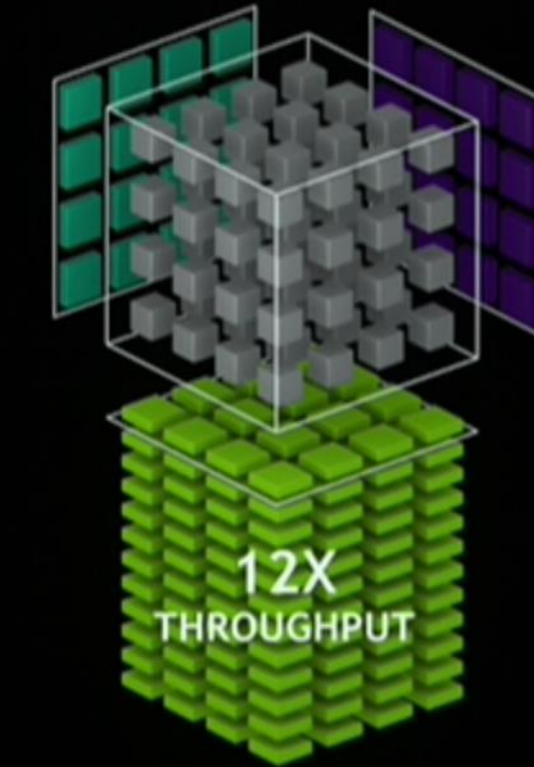# Deep Learning Everywhere



**Translation**



**Voice**



**Insurance**

# Heterogeneous Machine Learning Platform

**Application**

| E-commerce | Search | Speech recognition | Picture recognition | Big data | Mining Algorithm | Payment | Finance |

**Adaptive**

## Development Ecosystem

| IDE | Document | Sample | Community | Training | Operate & promote |

**Machine Learning Platform**

## Machine Learning PAAS

| Workflow (Arer) | model visualization | Authentication Authorization | Audit | Accounting | SDK |

## AIFLOW

## Offline training/Online Inference

| Open source Framework Tensorflow/MxNet/Caffe | In-house R&D Parameter Server |

| IO Optimize | RPC | MPI | Distribution |

## Heterogeneous computing IAAS

| GPU | FPGA | CPU | ASIC |

| CUDA Kernel extension | CUDA Kernel Optimize | Logic | Driver | CPU AVX | Cache Optimize |

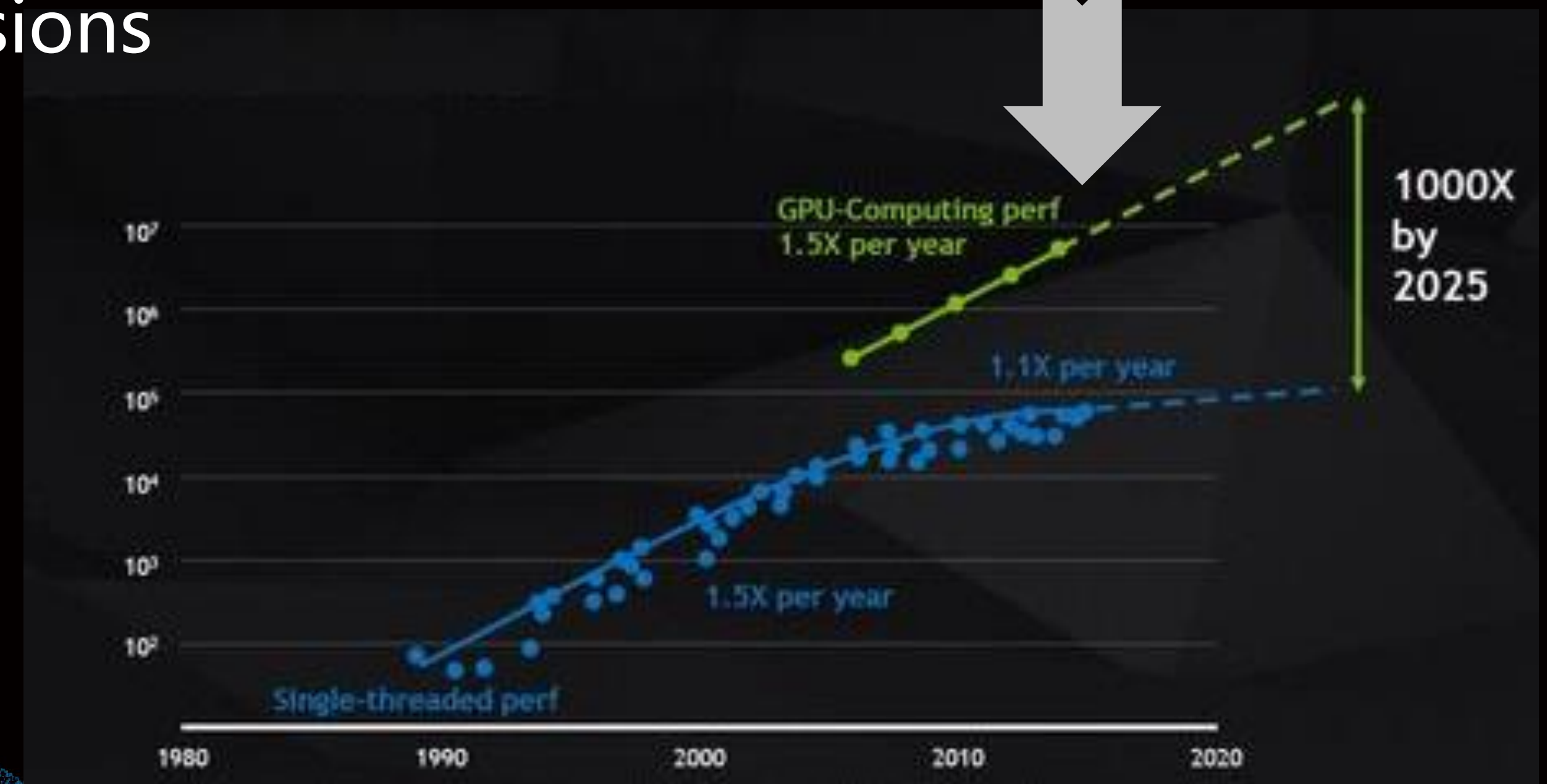| IDC | Networking | Virtualization（Docker） | Storage(ODPS、OSS) |

@ 2018 Alibaba Group

# Hardware Accelerated AI

- Training: Compute Intensive, Time Cost

- Inference: Service Oriented, Response Time

- Eco-System: Framework, Libs, Precisions
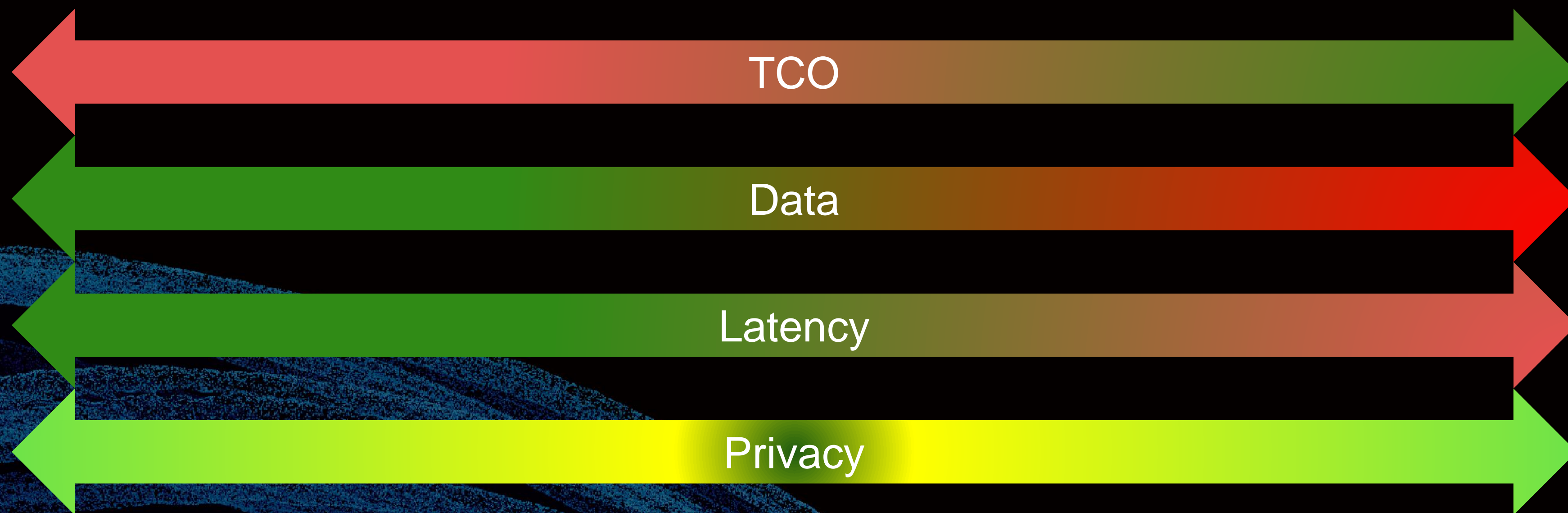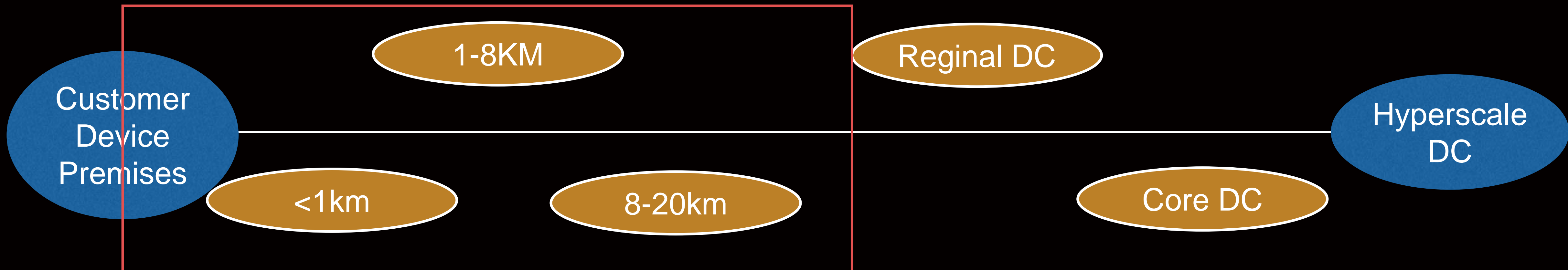
- Hardware Dividends for everyone

Tipping point:

- Google TPUs

- Volta TensorCore

- New hardware accelerators for AI
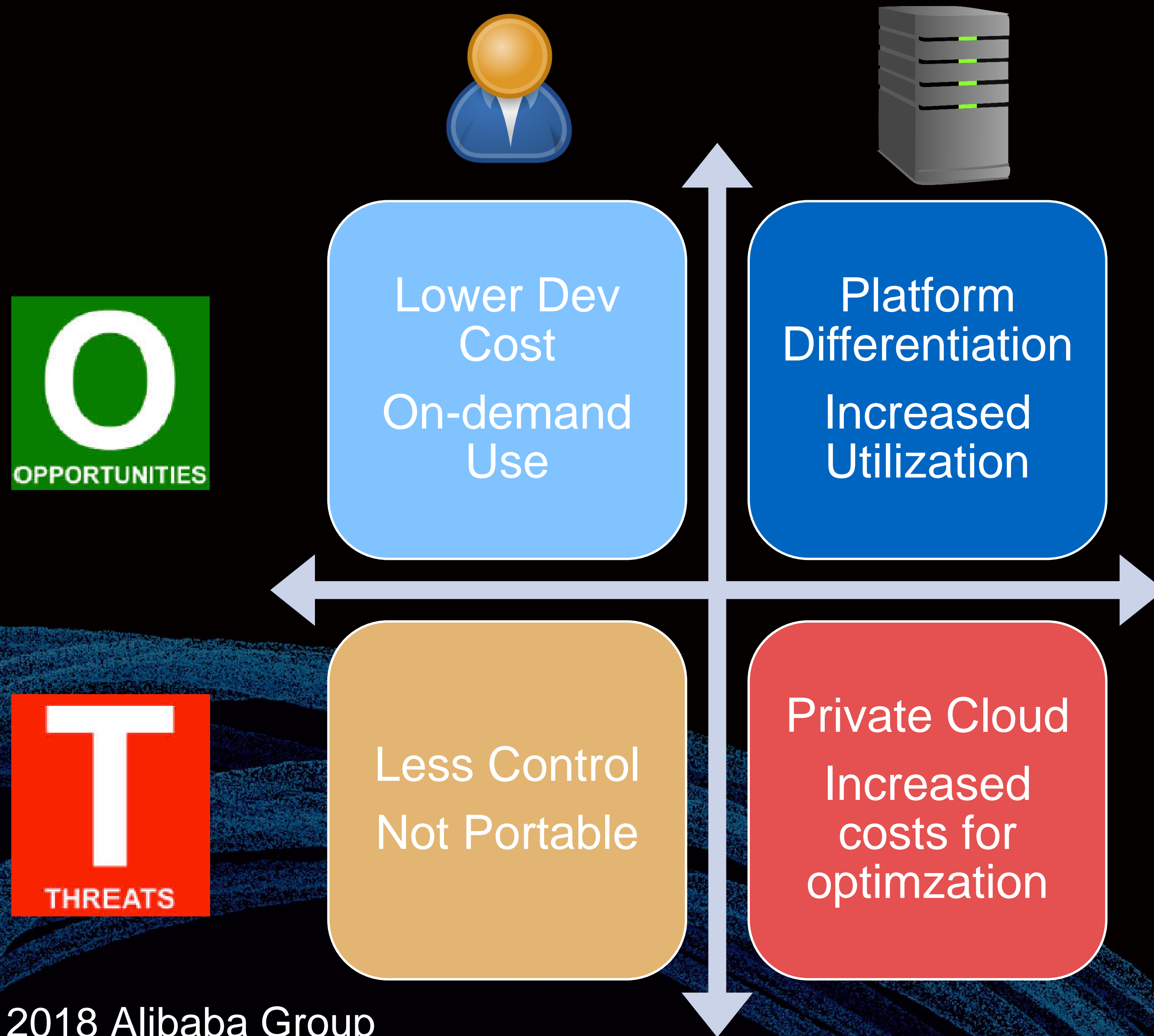
12X THROUGHPUT

GPU-Computing perf
1.5X per year

1000X
by
2025

$10^7$

$10^6$

1.1X per year

$10^5$

$10^4$

$10^3$

1.5X per year

$10^2$

Single-threaded perf

1980    1990    2000    2010    2020

*data from NVIDIA GTC 2017

# Edge – Forces of Gravity

Customer Device Premises

1-8KM

Reginal DC

<1km

8-20km

Core DC

Hyperscale DC

TCO

Data

Latency

Privacy

# Function Computing

**O**
OPPORTUNITIES

**T**
THREATS

| Lower Dev Cost<br>On-demand Use | Platform Differentiation<br>Increased Utilization |
| Less Control<br>Not Portable | Private Cloud<br>Increased costs for optimzation |

## Technology Challenges

- Quality of Service
- Infrastructure Utilization
- Accelerator Efficiency
- Capacity Granularity
- Multi-Tenancy Management
- Demand Projection
- Scheduling
- Compatibility

# Thank You!

E-mail: lingjie.xu@alibaba-inc.com