# 联想助力互联网智能转型

- 吴强 英伟达高级合作伙伴经理

# NVIDIA

> Founded in 1993

> Jensen Huang, Founder & CEO

> 13,227 employees

> $321B market cap;   **Q3 Rev:** $4.7B  YoY: 57%

"World's Most Admired Companies"
— Fortune

"50 Smartest Companies: #1"
— MIT Tech Review

"#1 Top CEO in the World"
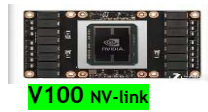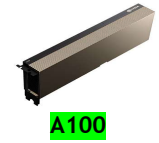— Harvard Business Review

"Most Innovative Companies"
— Fast Company

# TESLA / DGX 路线图



DGX-1　　　　DGX-Station　　　　DGX-2


**DGX-A100**
**DGX STATION A100**

---

K80
K40

K2
k1

M40
M4

M60
M10
M6

P100 nvlink
P100 PCI-e
P40
P4

**V100 PCI-e**

**Tesla T4**

**V100s**

**A100**

**V100 NV-link**

Kepler
Dynamic
Parallelism

Maxwell
DX12

Pascal
Unified Memory
3D Memory
NVLink 1.0

Volta
Unified Memory
Tensor Core
NVLink 2.0

Turing
RT Core +
Tensor Core

**RTX6000p**
**RTX8000p**

RT Core +
Tensor Core

**A40 PCI-e**
**RTX A6000**

2012　　　2014　　　2016　　　　2017　　　2018　　　2019　　　　2020

# GAME-CHANGING PERFORMANCE FOR INNOVATORS

## NVIDIA DGX A100 640GB System

10x Mellanox ConnectX-6 200 Gb/s Network Interface

500 GB/sec Peak Bi-directional Bandwidth

Dual 64-core AMD Rome CPUs and 2 TB RAM

3.2X More Cores to Power the Most Intensive AI Jobs

8x NVIDIA A100 GPUs with 640GB Total GPU Memory

12 NVLinks/GPU
600 GB/sec GPU-to-GPU Bi-directional Bandwidth

6x NVIDIA NVSwitches

4.8 TB/sec Bi-directional Bandwidth
2X More than Previous Generation NVSwitch

30TB Gen4 NVME SSD

50 GB/sec Peak Bandwidth
2X Faster than Gen3 NVME SSDs

### SYSTEM SPECIFICATIONS

| | NVIDIA DGX A100 640GB | NVIDIA DGX A100 320GB |
|---|---|---|
| GPUs | 8x NVIDIA A100 80 GB GPUs | 8x NVIDIA A100 40 GB GPUs |
| GPU Memory | 640 GB total | 320 GB total |
| Performance | 5 petaFLOPS AI 10 petaOPS INT8 | |
| NVIDIA NVSwitches | 6 | |
| System Power Usage | 6.5 kW max | |
| CPU | Dual AMD Rome 7742, 128 cores total, 2.25 GHz (base), 3.4 GHz (max boost) | |
| System Memory | 2 TB | 1 TB |
| Networking | 8x Single-Port Mellanox ConnectX-6 VPI 200Gb/s HDR InfiniBand 2x Dual-Port Mellanox ConnectX-6 VPI 10/25/50/100/200 Gb/s Ethernet | 8x Single-Port Mellanox ConnectX-6 VPI 200Gb/s HDR InfiniBand 1x Dual-Port Mellanox ConnectX-6 VPI 10/25/50/100/200 Gb/s Ethernet |
| Storage | OS: 2x 1.92 TB M.2 NVME drives Internal Storage: 30 TB (8x 3.84 TB) U.2 NVMe drives | OS: 2x 1.92TB M.2 NVME drives Internal Storage: 15 TB (4x 3.84 TB) U.2 NVMe drives |
| Software | Ubuntu Linux OS Also supports: Red Hat Enterprise Linux CentOS | |
| System Weight | 271.5 lbs (123.16 kgs) max | |
| Packaged System Weight | 359.7 lbs (163.16 kgs) max | |
| System Dimensions | Height: 10.4 in (264.0 mm) Width: 19.0 in (482.3 mm) max Length: 35.3 in (897.1 mm) max | |
| Operating Temperature Range | 5–30 ºC (41–86 ºF) | |

NVIDIA

# 双路CPU部分：AMD EPYC™ 7742
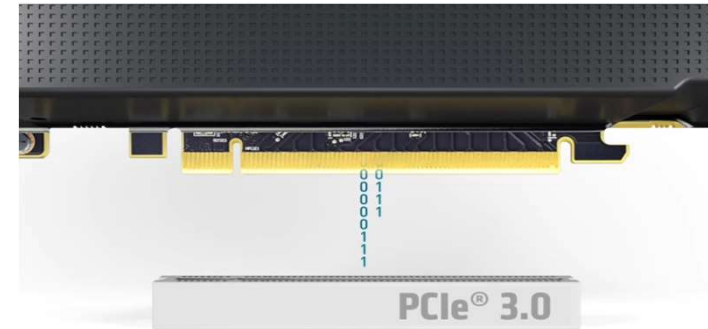
CPU 核心数量: 64
线程数量:128
基准时钟频率:2.25GHz

最大加速时钟频率:最高可达 3.4GHz
三级缓存:256MB
封装:SP3

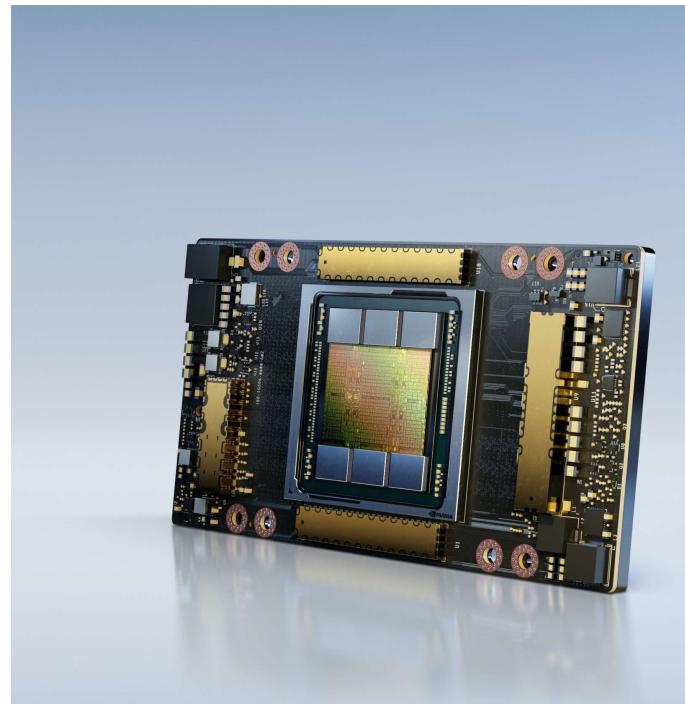支持的CPU插槽数:1P/2P
PCI Express 版本 :PCIe 4.0 x128
默认 TDP/TDP :225W

最高内存速度：Up to 3200MHz
内存类型：DDR4
内存通道：8
内存带宽（每路）：204.8 GB/s

PCIe® 3.0

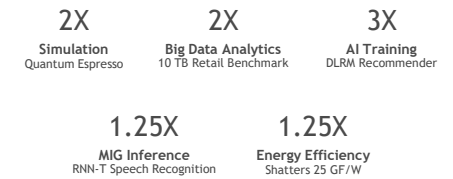## NVIDIA A100 GPU

Greatest Generational Leap – 20X Volta

| | Peak | | Vs Volta |
|---|---|---|---|
| FP32 TRAINING | 312 | TFLOPS TF32 | 20X |
| INT8 INFERENCE | 1,248 | TOPS | 20X |
| FP64 HPC | 19.5 | TFLOPS | 2.5X |
| MULTI INSTANCE GPU | | | 7X GPUs |

## SUPERCHARGED AI SUPERCOMPUTING WITH A100 80GB

World's Fastest GPU with World's Fastest Memory

A100 80GB Throughput vs A100 40GB

| 2X | 2X | 3X |
|---|---|---|
| Simulation | Big Data Analytics | AI Training |
| Quantum Espresso | 10 TB Retail Benchmark | DLRM Recommender |

| 1.25X | 1.25X |
|---|---|
| MIG Inference | Energy Efficiency |
| RNN-T Speech Recognition | Shatters 25 GF/W |

# 适用于办公室的服务器级解决方案

数据中心外的数据中心技术

- AI 时代的工作组服务器

- 没有数据中心但拥有数据中心性能

- 可以放置在任何地方的 AI 设备

  运行更大的模型，更快地获得答案

2.5 PFLOPS AI

320 GB GPU MEMORY

唯一具有 4 路 NVLink 和多实例
GPU（MIG）的工作组服务器

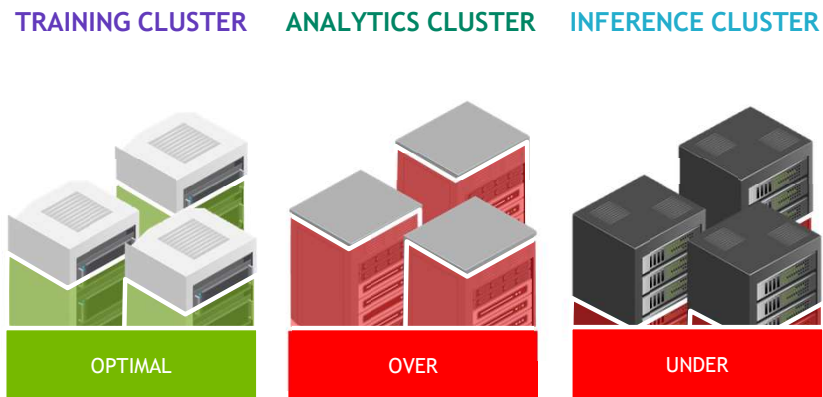| | DGX Station A100 320GB | DGX Station A100 160GB |
| --- | --- | --- |
| GPUs | 4x NVIDIA A100 Tensor Core GPUs | |
| GPU Memory (total) | 320GB | 160GB |
| Performance | 2.5 petaFLOPS AI; 5 petaOPS INT8 | |
| System Memory | 512GB DDR4 RDIMM, 3200MT/s | |
| Storage | OS: 1 x 1.92TB M.2 NVME<br>Data:1 x 7.68TB U.2 NVME | |
| CPU | AMD® Epyc® CPU 7742, 2.25GHz to 3.4GHz,<br>64 cores/128 threads, PCIe Gen4 | |
| Networking | Dual 10GBASE-T (RJ45) | |
| Display GPU | 4GB, 4x Mini DisplayPort | |
| Acoustics | <37dB | |
| Cooling | Custom refrigerant cooling system for GPUs and CPU | |
| System Power (max) | 1,5kW | |
| Management | AST2500, IPMI, Redfish | |
| System Dimensions | 518 D x 256 W x 639 H (mm) | |
| Operating Temp. | 5°C to 35°C (41°F to 95°F) | |

# ELASTIC AI INFRASTRUCTURE WITH DGX A100

DGX A100 with MIG Delivers New Agility for Today's Enterprise Data Center

**Traditional Infrastructure is Constrained**

Infrastructure silos starve AI workloads or waste capacity

**DGX A100 Infrastructure is Agile**

DGX A100 infrastructure uses MIG to allocate GPU resources to workloads
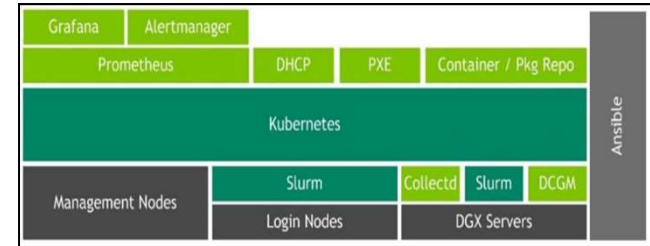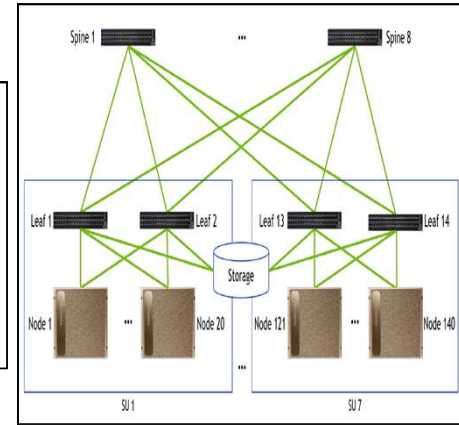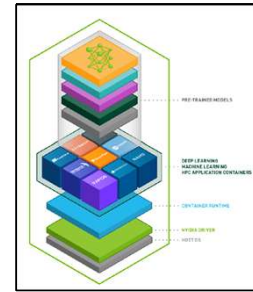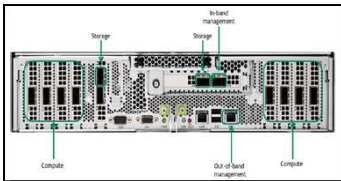
TRAINING CLUSTER    ANALYTICS CLUSTER    INFERENCE CLUSTER

OPTIMAL    OVER    UNDER

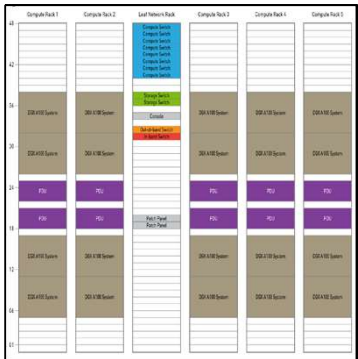Today    Tomorrow    Next Week

INFERENCE
TRAINING
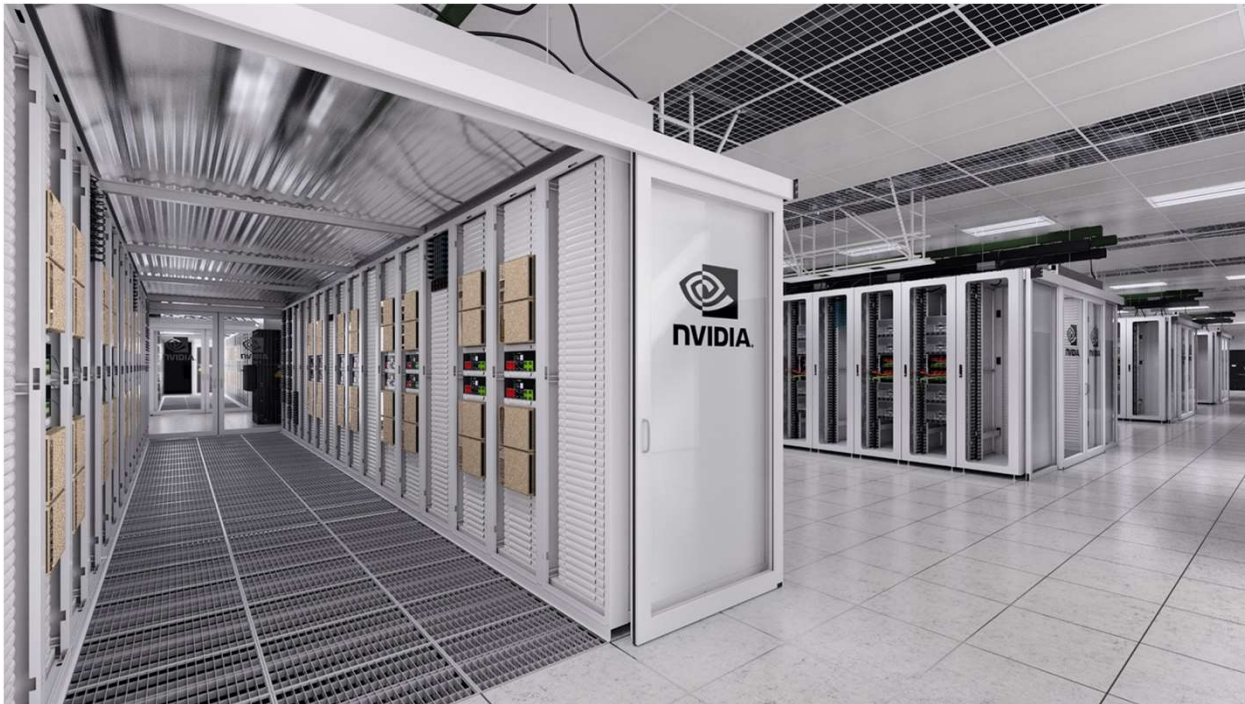ANALYTICS

# INTELLIGENTLY ADAPTED
# AND INTEGRATED

Flexible deployments tailored to the
unique needs of your environment

**DGX SuperPOD Solution for Enterprise:**

▸ Your partner to help your IT team navigate:

　　▸ AI workflow tools to speed time-to-insight

　　▸ Customized data center design, optimized for you

　　▸ Flexible deployment options tailored to your
　　　environment

▸ NVIDIA professional services in combination
　with
　select NVIDIA SuperPOD partners will
　speed your deployment experience

# DGX SUPERPOD
# DEPLOYMENTS AT NVIDIA

#1 on MLPerf for commercially available systems

#5 on TOP500 (63 PetaFLOPS HPL)

#1 on Green500 (26.2 GigaFLOPS/watt)
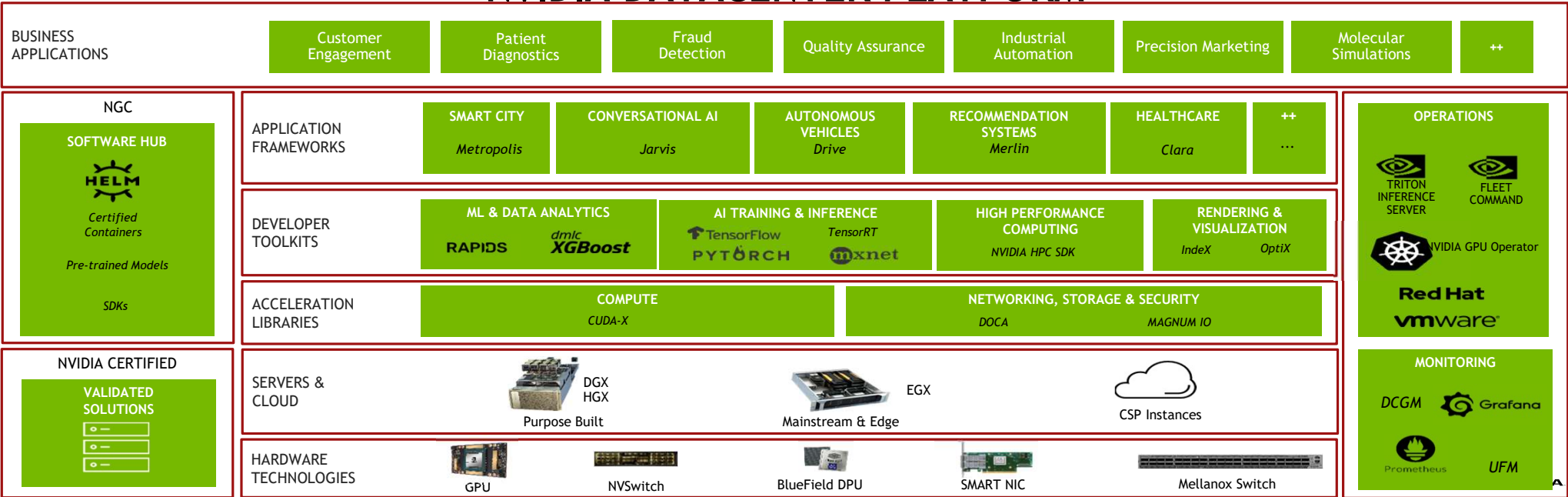
Fastest Industrial System in U.S.

Both are built with the NVIDIA DGX SuperPOD arch:

- ▶ NVIDIA DGX A100 and NVIDIA Mellanox IB
- ▶ NVIDIA's decade of AI experience

Selene Configuration:

- ▶ 4,480 NVIDIA A100 Tensor Core GPUs
- ▶ 560 NVIDIA DGX A100 640GB systems
- ▶ 850 Mellanox 200G HDR IB switches
- ▶ 14 PB of all-flash storage
- ▶ 2.8 ExaFLOPS of AI performance
- ▶ Built in 3 weeks

<span>⬢ NVIDIA.</span>

# NVIDIA DATACENTER PLATFORM

| BUSINESS APPLICATIONS | Customer Engagement | Patient Diagnostics | Fraud Detection | Quality Assurance | Industrial Automation | Precision Marketing | Molecular Simulations | ++ |
|---|---|---|---|---|---|---|---|---|

**NGC**

**SOFTWARE HUB**

HELM

*Certified Containers*

*Pre-trained Models*

*SDKs*

| APPLICATION FRAMEWORKS | SMART CITY **Metropolis** | CONVERSATIONAL AI **Jarvis** | AUTONOMOUS VEHICLES **Drive** | RECOMMENDATION SYSTEMS **Merlin** | HEALTHCARE **Clara** | ++ ... |
|---|---|---|---|---|---|---|

| DEVELOPER TOOLKITS | ML & DATA ANALYTICS RAPIDS *dmlc* XGBoost | AI TRAINING & INFERENCE TensorFlow TensorRT PYTORCH mxnet | HIGH PERFORMANCE COMPUTING NVIDIA HPC SDK | RENDERING & VISUALIZATION IndeX OptiX |
|---|---|---|---|---|

| ACCELERATION LIBRARIES | COMPUTE CUDA-X | NETWORKING, STORAGE & SECURITY DOCA MAGNUM IO |
|---|---|---|

**OPERATIONS**

TRITON INFERENCE SERVER
FLEET COMMAND

NVIDIA GPU Operator

**Red Hat**

**vmware**

**NVIDIA CERTIFIED**

**VALIDATED SOLUTIONS**

| SERVERS & CLOUD | DGX HGX — Purpose Built | EGX — Mainstream & Edge | CSP Instances |
|---|---|---|---|

| HARDWARE TECHNOLOGIES | GPU | NVSwitch | BlueField DPU | SMART NIC | Mellanox Switch |
|---|---|---|---|---|---|

**MONITORING**

*DCGM* Grafana

Prometheus *UFM*

# ACCELERATION SOFTWARE ON NGC

## 79+ Containers
DL, ML, HPC

## 27+ Model Training Scripts
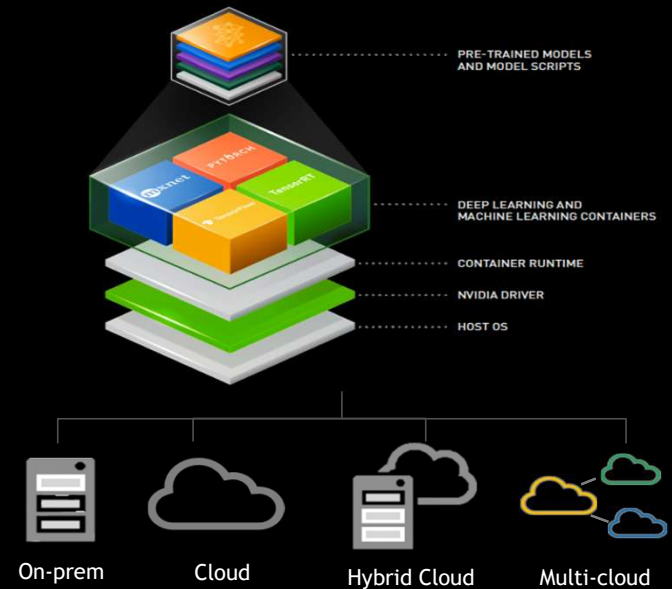NLP, Image Classification, Object Detection & more

NGC

## 60+ Pre-trained Models
NLP, Image Classification, Object Detection & more

## Industry Workflows
Medical Imaging, Intelligent Video Analytics

PRE-TRAINED MODELS AND MODEL SCRIPTS

DEEP LEARNING AND MACHINE LEARNING CONTAINERS

CONTAINER RUNTIME

NVIDIA DRIVER

HOST OS

On-prem     Cloud     Hybrid Cloud     Multi-cloud

# RECOMMENDERS — THE ENGINE OF THE INTERNET

**Billions of Users - Trillions of Items**

100's of millions of etail items -> Amazon & Alibaba recommenders

1000's of movies -> Netflix recommender

10's of millions of songs -> Spotify & iTunes recommenders

10's of millions of books -> Amazon recommender

Billions of Tik Tok & YT videos -> TT & YT recommender

Billions of websites -> Google search rank

So much news!!! -> Google & FB news recommenders

# KEY CHALLENGES IN CONSUMER INTERNET



## EXPLOSION OF AI MODELS & COMPLEXITY

- Number of models ranging from CNN, RNN, Transformer and new species like Wide & Deep

- Network complexity grew 10,000x in 7 years (peta flops)

## INCREASING DATA ONSLAUGHT

- Datasets continuing to increasing dramatically

- Multiple sources, different formats, varying quality

## CUMBERSOME DEPLOY & MANAGE WORKFLOWS

- Scaling infrastructure monitoring & management is a challenge

- Software deployment, versioning and updates are time consuming

# GPU在互联网行业的主要应用领域

| Content Understanding | Conversational AI | Recommender and Search | Video Processing | Machine/Drive | Data Analysis |
|---|---|---|---|---|---|
| • 内容审核<br>• 内容标签/分类检测<br>• 语音输入<br>• 语音识别<br>• 文本翻译<br>• 人脸/物体检测/识别<br>• OCR | • 智能终端<br>• 智能客服<br>• 在线会议/社交<br>• 语音助手<br>• 各类语音合成<br>• 智能问答<br>• 虚拟主播 | • 推荐/搜索/广告 CTR<br>• 召回、排序<br>• 内容标签挖掘和向量化 | • 编解码<br>• 渲染特效<br>• 质量增强<br>• 视频编辑<br>• 生成特效<br>• AR/VR<br>• Cloud Gaming | • 自动驾驶<br>• 智能配送<br>• 物流机器人<br>• 智能货架<br>• 仿真模拟 | • Spark<br>• 数据预处理<br>• 特征提取<br>• 传统机器学习<br>• 图计算<br>• 数据可视化 |
| CV, ASR, NLP | ASR，TTS，NLP | CTR、NLP, CV | Codec、CG、GAN、CV | Jetson、CV、speech、CG | Bigdata/Spark、xgBoost/ML |

## 主要类别互联网公司业务分布

| | Content understanding | Conversational AI | Recommender and search | Video Processing | Machine/Drive | Data Analysis |
|---|---|---|---|---|---|---|
| 云服务2B | ＊＊＊＊ | ＊＊＊ | | ＊＊ | | ＊＊ |
| 电商 | ＊＊＊＊ | ＊＊ | ＊＊＊＊＊ | ＊ | ＊＊＊ | ＊＊ |
| 社交/效率办公 | ＊＊＊ | ＊＊＊ | ＊＊ | ＊＊ | | ＊ |
| 娱乐/视频 | ＊＊＊＊＊ | ＊ | ＊＊＊＊＊ | ＊＊＊＊ | | ＊ |
| 教育/音频 | ＊＊＊ | ＊＊ | ＊ | ＊＊＊ | | ＊ |
| 生活服务 | ＊＊＊＊ | ＊＊ | ＊＊ | ＊ | ＊＊＊ | ＊＊ |
| 资讯/搜索 | ＊＊＊＊ | ＊＊ | ＊＊＊＊＊ | ＊＊ | | ＊＊ |
| 金融 | ＊＊ | ＊ | ＊ | | | ＊＊＊ |
| 手机 | ＊＊＊＊ | ＊＊ | ＊＊ | ＊＊ | | ＊ |

**REAL-TIME SPEECH SERVICES AT SCALE**

WeChat, a leading Chinese social media platform with ~1B users, wanted to improve its speech to context services. But as the company deployed its new acoustic model, its CPU-only servers were unable to effectively run the new version. WeChat deployed servers equipped with Tesla P4 GPU inference accelerators and increased speech inference throughput by 2.5X and in-model accuracy by 20% — all while staying within its low latency budget.

推荐分类

美女热舞 HOT

绝地求生 HOT

娱乐

| 音乐 | 舞蹈 |
| 脱口秀 | 户外 |
| 搭讪 | 喊麦 |
| 二次元 | 体育 |
| 美食 | 展开全部 ˅ |

游戏

| 绝地求生 | 刺激战场 |
| 王者荣耀 | 主机热游 |

综合

安卓下载YY极速版
登录立拆7元红包

手机YY    PC YY

不想撞南墙 想撞你心上
QS、倩倩                15.8万

大鹏说事
舞帝大鹏              6501

子航喜乐汇
舞帝子航              7567

华矩公会风云排行榜
华矩阿狼              1142

有时间，来听听别人的...
海郎中              15.2万

心悦小忆送欢乐
6961小忆              1229

小贺舞帝榜单
舞帝秀子贺              12141

温婉秀美以南酱
音豪@小以南（慧慧...）              131

云源-源徒小纯皇亲哦戚
雲源- 演纯源徒 皇亲...              10782

又出什么大事了劲··爆··
白主任              1339

周播大人物

新青铜小可可
3可可              434

## SAFEGUARDING LIVE STREAMING CONTENT

YY offers 100M concurrent live stream participants fun and engaging experiences. Auditing content during live streams to detect and filter inappropriate material requires real-time inference.

By deploying NVIDIA TensorRT on GPUs, YY achieved 30% higher inference workload throughput and reduced memory requirements by 40%.

Real-time inference prevents inappropriate content from making its way into live streams.

NVIDIA    YY.COM

## THE 'DATA' SCIENCE OF PERSONALIZED SEARCH

Pinterest, one of the largest global social media companies, allows users to discover information using images, GIFs and video. The company continually works to improve personalized search results for its 250M monthly users.

To leverage AI for this Pinterest processes data from users and 175+B pins stored on AWS. GPU-powered deep learning on AWS EC2 P3 speeds training from months to days and allows Pinterest to quickly scale AI projects. The outcome is personalized search results for 600+M visual searches per month.

**REAL-TIME FRAUD DETECTION**

When PayPal was looking to deploy a new fraud detection system, they set a high bar: the system had to operate worldwide 24/7 and in real-time to protect customer transactions from potential fraud

CPU-only servers couldn't meet the requirements

Using NVIDIA T4 GPUs, PayPal delivered a new level of service, using GPU inference to improve real-time fraud detection by 10% while lowering server capacity by nearly 8X

# 云游戏行业应用场景