The Practices and Learning
on AI Edge Computing

Gang Chen | Network System Architect
25th June, 2019
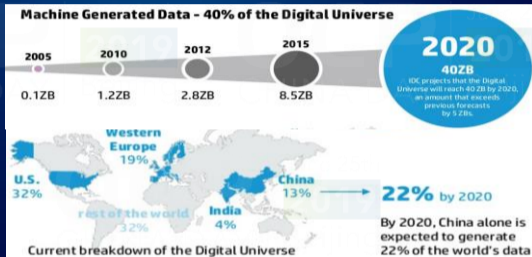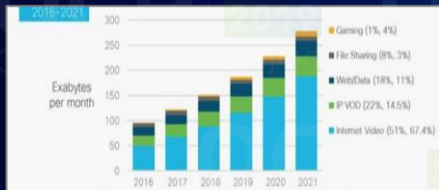
**CONTENTS**

# Streaming and Computation Schema

Computation: most data will be computed on Edges

Streaming: growing steadily at 45% per year





Machine Generated Data – 40% of the Digital Universe

Definition : 1080P -> 4K -> 8K

Fluency : 30fps -> 60fps ->120fps

Multi-Flows: Single Flows -> 360° immersive exp

Model: on-demand -> Live Broadcast

Industrial devices equipped with edge computing board

Smart Phone

→ Edge Computing@5G & CDN

# Consumer && Industrial Internet

## Consumer Internet Evolution

**Streaming consumption upgrading:**
*UHD on-demand, no cache, streaming socialization*

**Scenario-oriented AI population:**
*AR trick, picture and video rendering, editing*

**Multi-mode interaction:**
*Wearable devices, vehicles entertainment, AR/VR glass*

**Trend I : Computation on D-E-C**
*Computation load balance on Device- Edge- Clouds*

**Trend II : Multi-mode services**
*Emphasis on "Last Mile" technology and devices*

## Industrial Internet Revolution

**Industrial networking:**
*Requirements on algorithm, computation and intelligent*

**Infra sharing & opening platform:**
*From vertical Silo to digital platform services*

**Digital & physical world integration:**
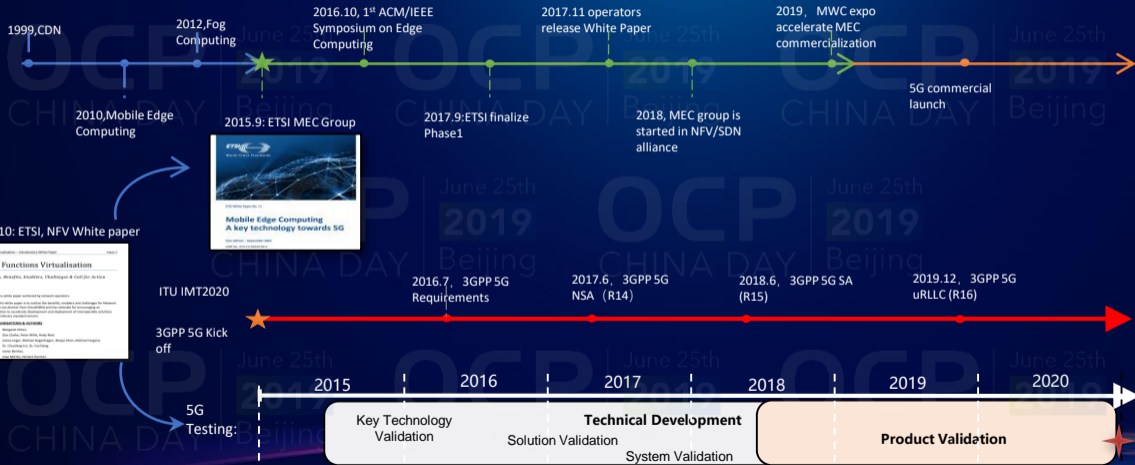*Smart city, Smart transport, AI security, etc,.*

**Trend III : The middle-ware platform for industry**
*Open Source for de facto standard, Industrial OS*

**Trend IV : Networking & Flattening**
*Data growth triggers computation flattening*

# Edge Computing Powered by BAIDU AI Full-stack

**AI Solutions** 解决方案

| 人脸会场签到 | 人脸身份验证 | 人脸会员识别 | 人脸闸机 | 智能安防视频分析 | 机器人视觉 | 文本内容审核 | 智能电销 | 视频内容审核 | 图像/头像审核 | 视频内容分析 |

**AI Gears** 软硬一体

| 远场语音开发套件 | 人脸识别开发套件 | 机器人导航与视觉模组 | 服务机器人 |

**AI Capabilities** 场景化能力

云服务 / 离线SDK / 服务器私有部署

| 语音 | 人脸 | 文字识别 | 图像识别 | 内容审核 | 视频技术 | 自然语言处理 |
|---|---|---|---|---|---|---|
| 实时语音识别 | 手势识别 | 出租车票识别 | 通用文字识别(高精度) | 图像超分辨率 | 图像对比度增强 | 定制化评论点抽取 | 文本审核 |
| 音频文件转写 | 驾驶行为分析 | 火车票OCR | | 食材识别 | 图像去雾 | 定制化词法分析 | 汉语纠错 |
| 医疗语音识别 | 人像分割 | 数字OCR | 通用文字识别(含生僻字) | 相册分类 | 水印/二维码/条形码提取 | 词义相似度 | 文本审核 |
| 呼叫中心识别 | 人流量统计 | 彩票OCR | | 水印识别 | | 依存句法分析 | 对话情绪识别 |
| 远场语音识别 | 人体关键点 | 表格OCR | 通用文字识别(含位置) | 商品图像搜索 | 公众人物识别 | DNN语言模型 | 文章标签 |
| 长音频识别 | 人体属性 | 银行卡识别 | 网络图片文字识别 | 相同图像搜索 | 图像质量检测 | | 文章分类 |
| 手机话筒识别 | 活体检测 | 身份证识别 | 名片OCR | 相似图像搜索 | 恶心识别 | 词向量表示 | 评论观点抽取 |
| 语音唤醒 | 人脸M:N | 护照识别 | 手写OCR | 场景与物体识别 | 政治敏感识别 | | 情感倾向分析 |
| 语音合成 | 人脸1:N | 车牌OCR | 增值税发票OCR | 主体检测 | 色情识别(GIF) | 词向量表示 | 短文本相似度 |
| | 人脸1:1 | 营业执照OCR | 二维码识别 | 通用物体分类 | 暴恐识别 | 词法分析 | 文章分类 |
| | 人脸属性 | 行驶证 | 票据OCR | 动物识别 | 色情识别 | | |
| | 人脸检测 | | 驾驶证 | 花卉识别 | | | |
| | | | | 植物识别 | | | |
| | | | | 车型识别 | | | |
| | | | | 菜品识别 | | | |
| | | | | Logo识别 | | | |

第三方服务

**AI Platform** 平台

| 语音技术 | 人脸与人体识别 | OCR | 图像识别 | 图像审核 | 视频技术 | 自然语言处理 |

| EasyDL—定制化模型训练与服务平台 | UNIT—智能对话系统开发平台 | 自定义模板文字识别平台 | 机器翻译开放平台 | AR/VR开放平台 | 数据智能平台 |

**AI Framework** 框架

| 教育 / 生态开源 | PaddlePaddle 训练营 | 飞桨开放 | 室外场景理解 | 视频精彩片段 | 阅读理解 | 信息抽取 | 知识抽取 | 交通预测 |
| | 深度学习框架 PaddlePaddle | | | | | | |

# Computation Deployment: DEC Model

Device ⟷ Edge ⟷ Cloud

## D-E-C Computation Model

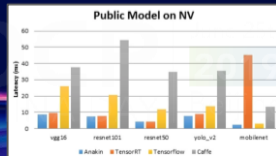| | Device | Edge | Cloud |
|---|---|---|---|
| 算力 | 1-10 TOPS | 10-100 TOPS | 100+ TOPS |
| 功耗 | 0.1-10 W | 10-100 W | 100+ W |
| 延时 | 10-100 ms | ms~s | ms~s |

AI HW Accelerators

Up to 260Tops



AI IF Accelerators

Support Intel-CPU, NV-GPU, AMD-GPU and etc,.
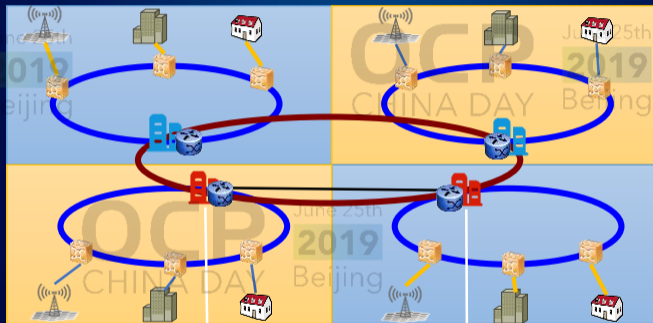
# Network Deployment: Grid Model



Campus, 60k- 70k
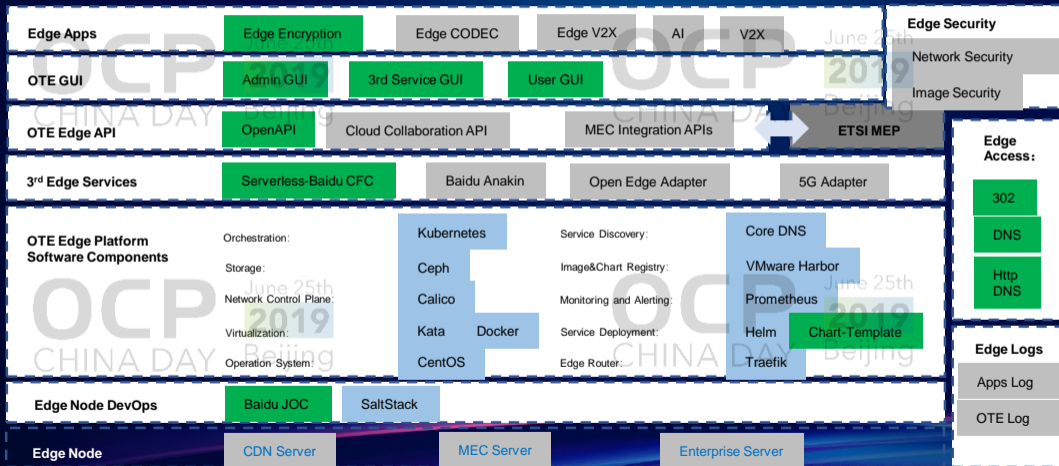
Metro-access, 5k-7k

Metro-core, 600-700

Provincial , 20-30

# OTE: Over The Edge

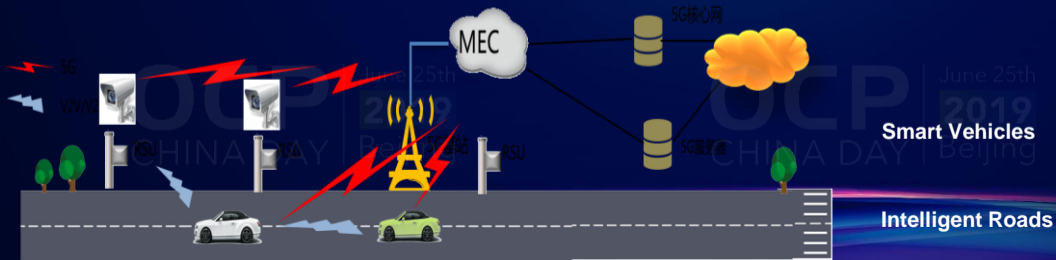| Edge Apps | Edge Encryption | Edge CODEC | Edge V2X | AI | V2X |
|---|---|---|---|---|---|

**Edge Security**
- Network Security
- Image Security

| OTE GUI | Admin GUI | 3rd Service GUI | User GUI |
|---|---|---|---|

| OTE Edge API | OpenAPI | Cloud Collaboration API | MEC Integration APIs | ETSI MEP |
|---|---|---|---|---|

| 3rd Edge Services | Serverless-Baidu CFC | Baidu Anakin | Open Edge Adapter | 5G Adapter |
|---|---|---|---|---|

**OTE Edge Platform Software Components**

| | | | |
|---|---|---|---|
| Orchestration: | Kubernetes | Service Discovery: | Core DNS |
| Storage: | Ceph | Image&Chart Registry: | VMware Harbor |
| Network Control Plane: | Calico | Monitoring and Alerting: | Prometheus |
| Virtualization: | Kata    Docker | Service Deployment: | Helm    Chart-Template |
| Operation System: | CentOS | Edge Router: | Traefik |

**Edge Access:**
- 302
- DNS
- Http DNS

**Edge Logs**
- Apps Log
- OTE Log

| Edge Node DevOps | Baidu JOC | SaltStack |
|---|---|---|

| Edge Node | CDN Server | MEC Server | Enterprise Server |
|---|---|---|---|

CONTENTS

# V2X @ MEC

| C-V2X | Access | Bandwidth | Latency | Request MEC | Mobility |
|---|---|---|---|---|---|
| perception | 5G/Fiber | 8Mbps（上行） | 5ms | Y | N |
| RSU | 5G/Fiber | 1Mbps（下行） | 5ms | Y | N |
| Automotive | 5G | 50-100Mbps | 10ms | Y | Y |
| Vehicle Entertainment | 5G | 20-50Mbps | 20~30ms | Y | Y |

# UHD 8K Live Broadcast @ MEC

- ## 8K UHD Broadcast based on 5G SA



- >100Mbps symmetric streaming

- Live AI analytic processing

- Multiply streaming protocols supports

# AR Rendering @ MEC

- Online Translations and AR rendering are completed on the edges

| | 4G | 5G |
|---|---|---|
| Delay | 24ms | 5.8~23ms |
| BW | 2.32Mbps | 600Mbps |

OCR

AR render

Streaming Analytics

MEC
| AI | Fast Path | Storage Acceleration | Codec | Encryption |

Virtualization Environment

Networking

Computing

Storage

Rendering Stream

5G

# Thank you