



北京
2019

遨游“视”界 做你所想
Explore World, Do What You Want

LiveVideoStackCon 2019 北京

2019.8.23-24

出品: LiveVideoStack CSDN
—— 音视频技术社区 ——

xNN: 支付宝App中的实时AI引擎

蚂蚁金服 - 多媒体产品技术部 - 机器视觉
周大江 (弘川)



深圳
2019

遨游“视”界 做你所想
Explore World, Do What You Want

LiveVideoStackCon 2019 深圳

2019.12.13-14



出品: **LiveVideoStack**
—— 音视频技术社区 ——

成为讲师: speaker@livevideostack.com

成为志愿者: volunteer@livevideostack.com

赞助、商务合作: kathy@livevideostack.com

端侧机器学习

机遇

端侧固有
资源限制



支付宝App
场景难点

挑战

体验

使能**实时**交互应用，
体验不受网络影响，
部分场景端到端精度
更高

成本

计算本地化，节省云
侧计算/流量/存储开
销

隐私

数据上传非必须，保
护用户敏感数据，规
避商户法律与舆论风
险

算力

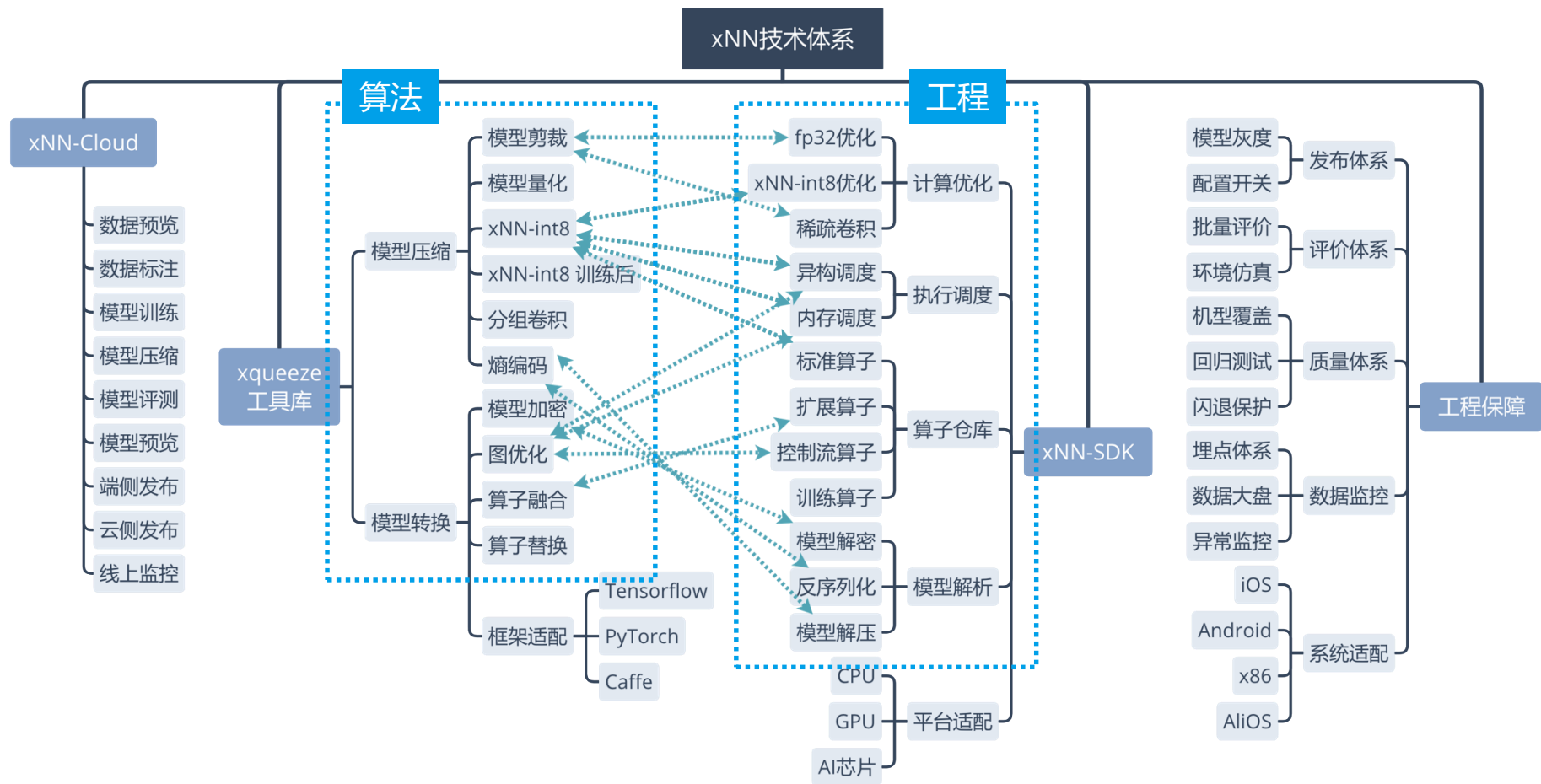
海量用户基数与“普
惠”服务原则，更需
要对**低端老旧机型**的
不离不弃

空间

支付宝App承载众多
业务线，空间资源的
紧张程度远超一般应
用

效率

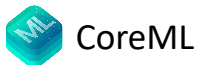
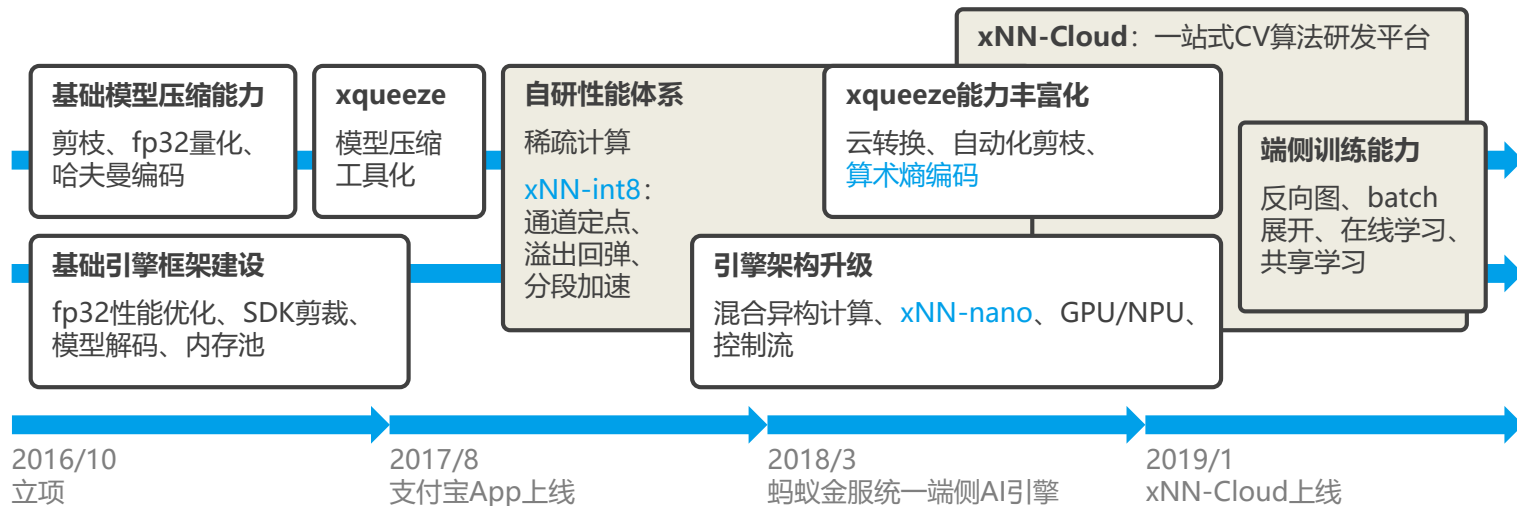
如何以有限的人力支
撑众多业务线的端侧
机器学习需求



xNN技术演进

算法

工程



xNN模型压缩

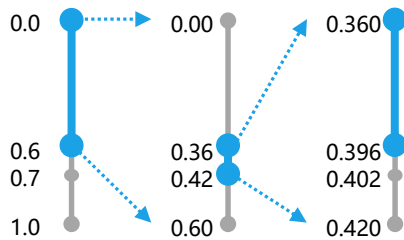
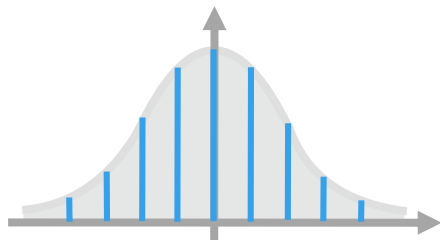
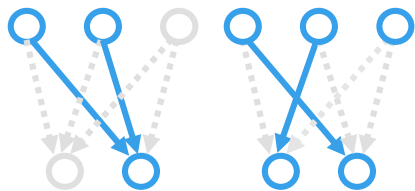
算法仓库
构成

神经元/突触剪枝
模型结构精简化

xNN-int8体系
模型参数离散化

自适应算术编码
逼近理论压缩极限

较Huffman提升15%



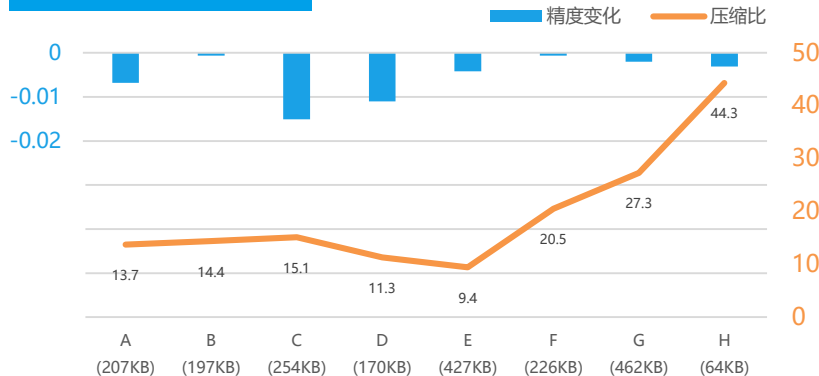
算法使用体验

```
import tensorflow as tf

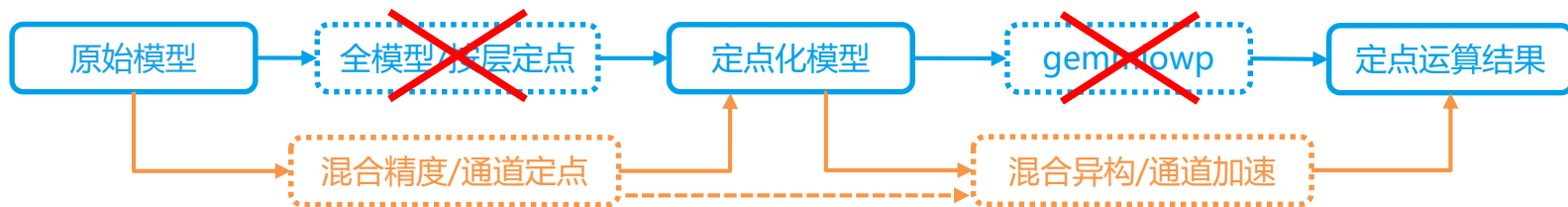
def train():
    train_op, loss = mobilenet_v1()
    with tf.Session() as sess:
        if step < 100000:
            step += 1
            _, val = sess.run([train_op, loss])

if __name__ == '__main__':
    train()
```

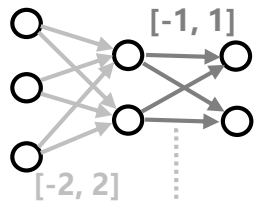
模型压缩效果



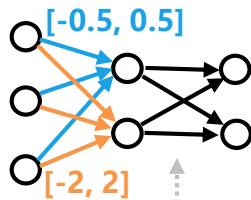
xNN-int8体系：精度提升



逐层内参数统一定点



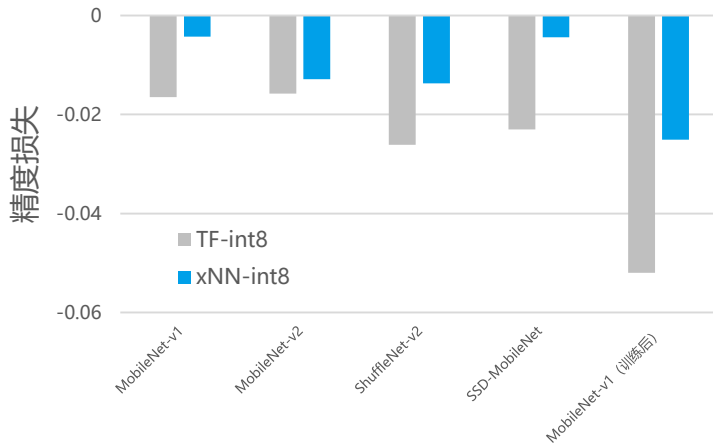
按输出通道定点



回归层精度敏感，允许
int8/float32混合精度

通道间分布不均
常见于移动端模型结构如：
Group Convolution
Depth-wise Convolution

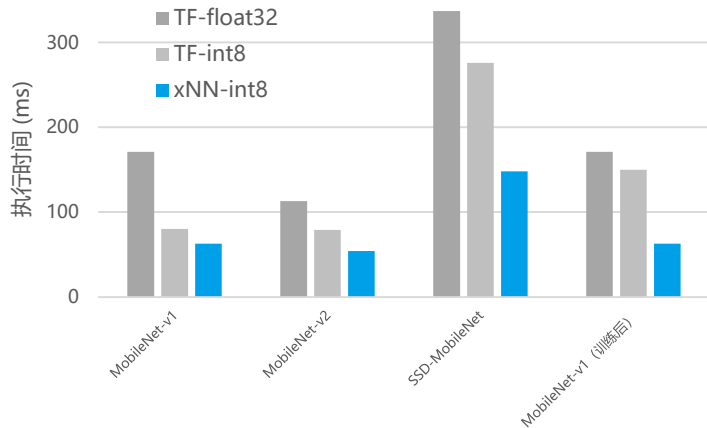
xNN-int8精度普遍提升~1%



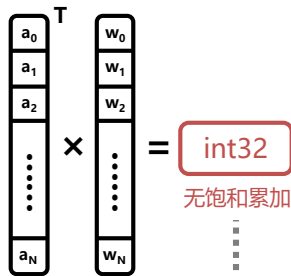
xNN-int8体系：性能提升



xNN-int8性能：提升~20%

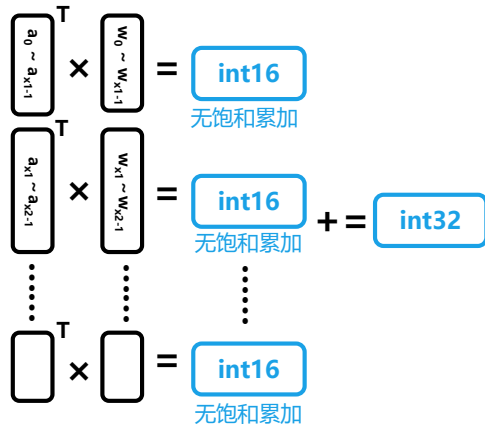


传统定点实现



int32较int16累加器
多50+%指令开销

幅值分段-无饱和加速



与业界方案的对比

计算性能

模型压缩

SDK尺寸

工具链

xNN方案

定制int8体系
训练中+训练后量化
精度高+速度快

剪枝+量化+算术熵编码
10-30倍压缩比

去STL依赖
最简配置不到100KB

标注-训练-压缩-转换-测试-
部署-监控 全链路可视化

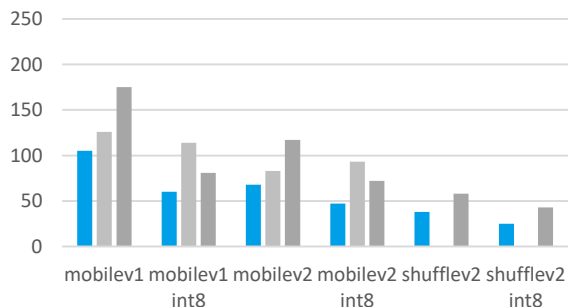
业界主流

标准int8体系
训练后量化

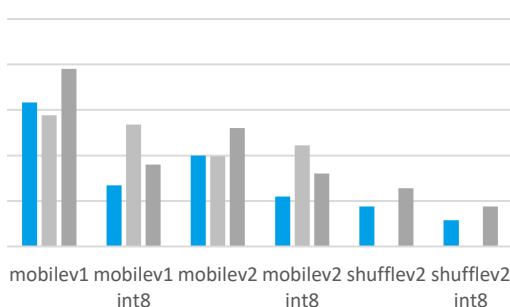
仅量化
2-4倍压缩比

300KB+

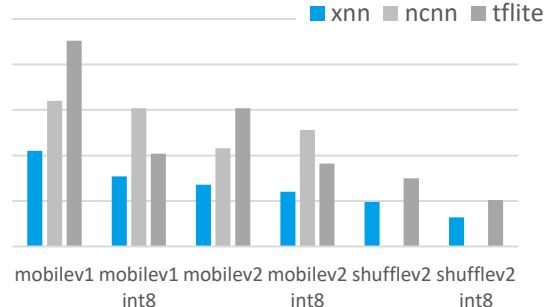
命令行转换



cpu_time@Mi6 (ms)



cpu_time@vivo X20A (ms)

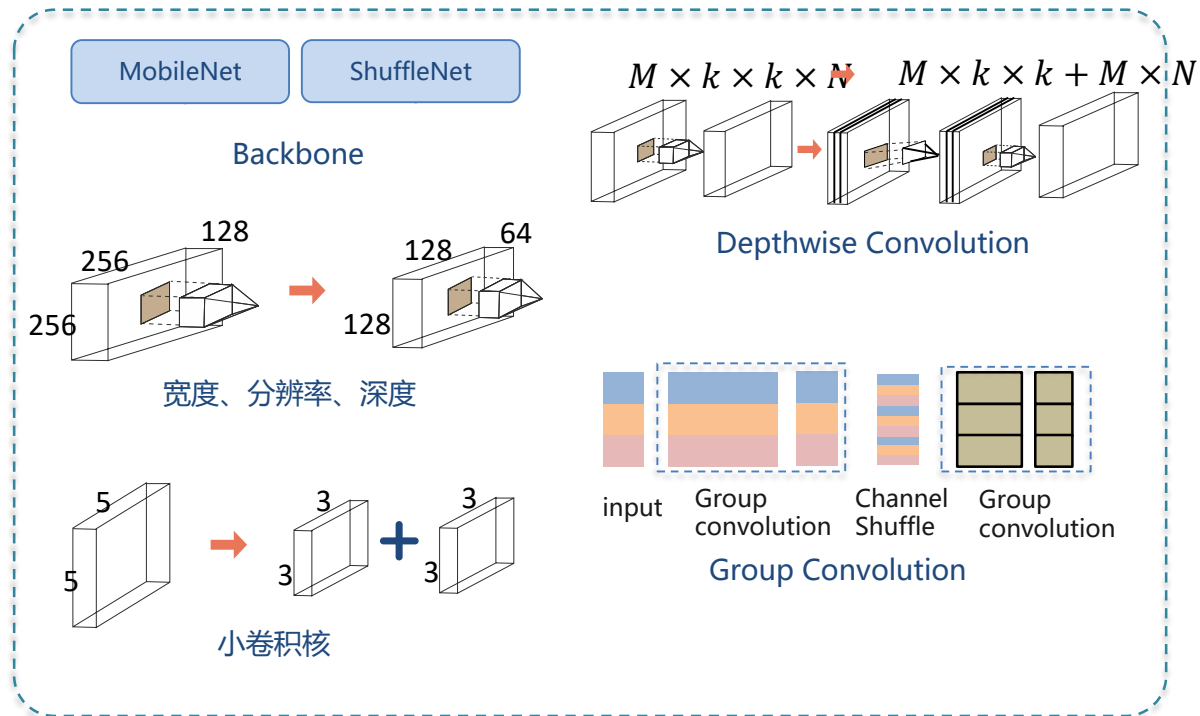


cpu_time@vivo X23 (ms)

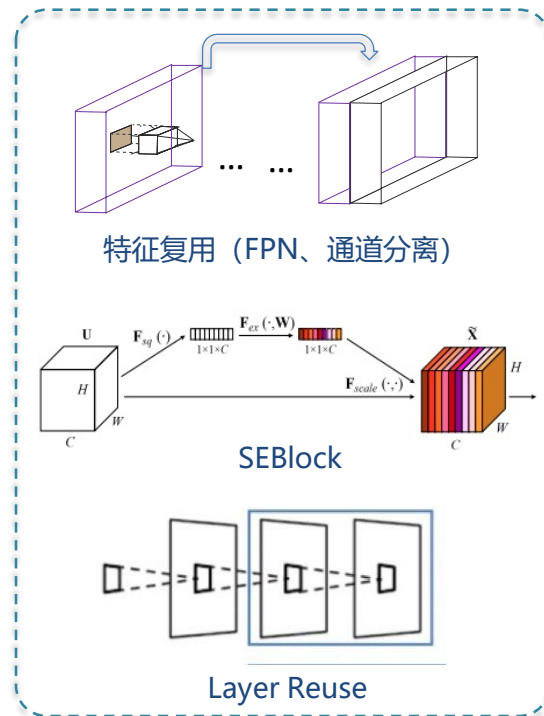


高效/轻量级网络设计

- 高效网络设计思想



- 轻量化和精度平衡



端侧实时AI能力：金融信息场景



将银行卡卡面放在此区域，扫描卡片

银行卡OCR

识别准确率 99+%
模型尺寸 460KB



身份证OCR

6000 字中文字库
全字段准确率 84+%
模型尺寸 980KB



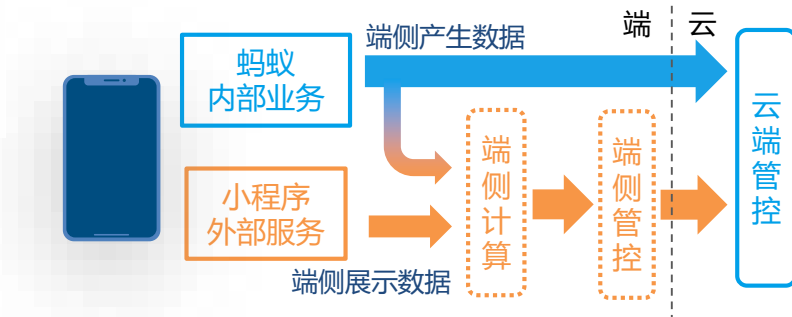
理赔宝

实时拍摄引导
图片质量
文档属性/类型
有效数据比例：
40% → 90%



定损宝

车牌/VIN码OCR
距离检测
模糊/光线检测
部件/损伤检测



端+云 内容安全

通用OCR 1.2MB
端侧分流50%计算

鉴黄模型 253KB
端侧减少70%数据上报



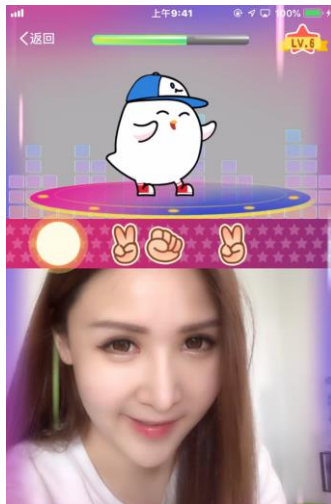
端侧实时AI能力：互动运营场景



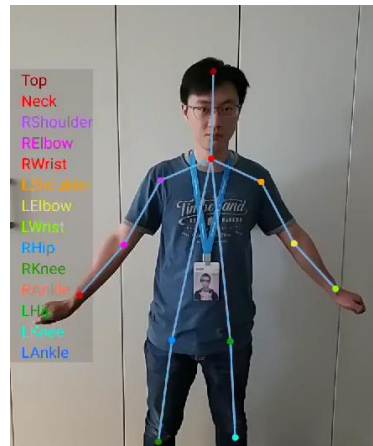
春节扫福



AR会议室



庄园手保健操



人体姿态



手势关节点

xNN应用

40+ 业务场景、80+ 模型 线上运行

支付宝App



新春红包：2019年活动期间支持3000亿+调用，客户端覆盖率98%



扫码支付：码姿态矫正+小码放大，显著提升扫码体验

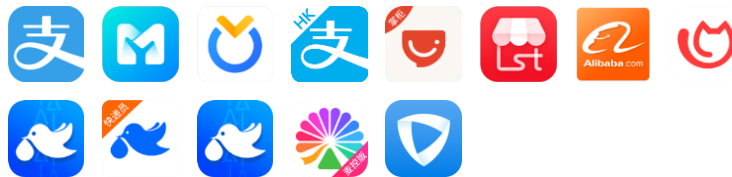


智能加载：支付宝App卡顿减少，关键业务启动速度提升



资金风控：通过手机操作预测丢失，保障用户资金安全

阿里集团App

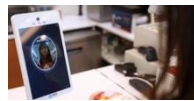


合作银行&保司：10+ App

IoT设备场景



刷脸支付：使用xNN-int8后性能翻倍，显著提升刷脸体验

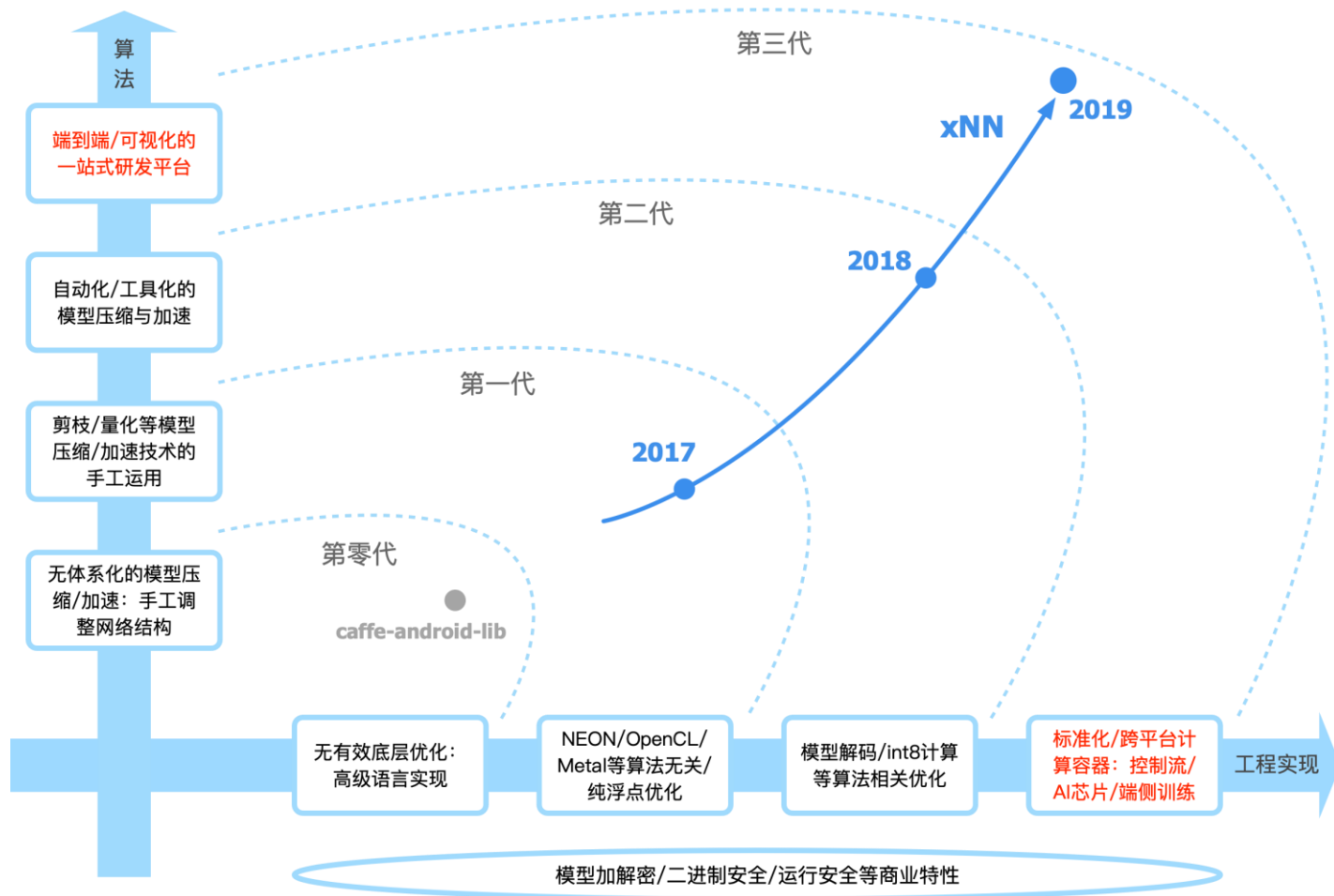


“蜻蜓”：刷脸支付+扫码支付，体验同步提升



高德车镜：10万台设备采用xNN实现端侧计算

端侧AI的演进



xNN-Cloud: 一站式CV算法研发平台



数据标注

人工标注
模型标注

模型训练

调度控制
参数搜索

端化处理

压缩加速
模型加密

模型表现

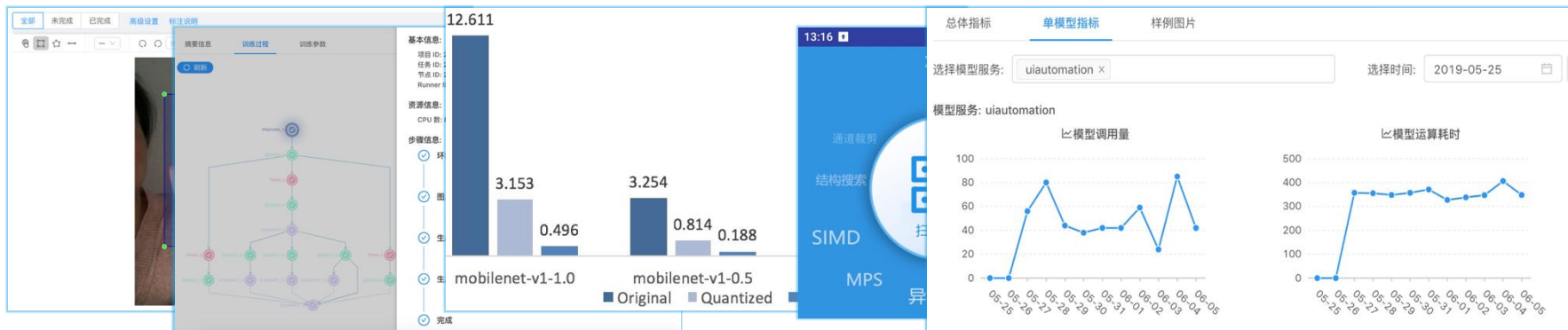
手机预览
云侧平台

模型发布

客户端
服务端

线上监控

性能指标
数据回流



XNN-CLOUD

5分钟带您了解

[XNN.ALIPAY.COM](https://xnn.alipay.com)

xNN-Cloud: 算法能力沉淀

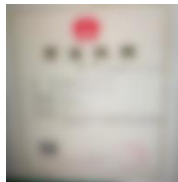
- 平台化沉淀算法模型，用户可自助训练业务模型

分类



识花

大小: 253 K
精度: 0.98
耗时: 73 ms



营业执照

大小: 231 K
精度: 0.99
耗时: 62 ms



身份证

大小: 226 K
精度: 0.99
耗时: 57 ms

...

检测



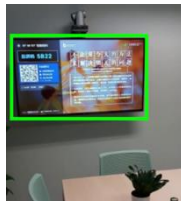
Landmark

大小: 312 K
精度: 0.95
耗时: 125 ms



Icon

大小: 220K
精度: 0.97
耗时: 115 ms



屏幕

大小: 230 K
精度: 0.99
耗时: 120 ms

...

OCR



文档标题

大小: 628 K
精度: 0.9
耗时: 320 ms



快递单

大小: 980 K
精度: 0.85
耗时: 600 ms



投屏码

大小: 530 K
精度: 0.95
耗时: 240 ms

...

公开参考资料

这款神秘的移动端OCR引擎，如何做到“所见即所得”？

原创：亦弦 阿里技术 4月2日



阿里妹导读：随着深度学习，尤其是CNN和RNN等技术的飞速发展，文字识别技术(OCR)识别精度越来越高。与此同时，在保护个人隐私和隐私安全越来越重要的今天，OCR也面临着新的挑战。本文介绍支付宝App中的深度学习引擎——xNN-OCR。



不写一行代码，完成机器视觉算法的研发

原创：南烽 阿里机器智能 3月26日



小叭导读：xNN云平台面向普通的研发、质量、UE等角色，也面向视觉、语音等提供全自动的模式。配合专业的数据标注工具、端侧发布、部署、运维能力，使得研发、测试、发布、运维全流程的研发流程。



含代码 | 支付宝如何优化移动端深度学习引擎？

原创：家大 阿里技术 2018-06-20



阿里妹导读：移动端深度学习在增强体验实时性、降低低端计算负载、保护用户隐私等方面具有天然优势。在支付、安全等领域，考虑到移动端落地面临的性能、内存使用、模型大小等挑战，本文介绍支付宝App中的深度学习引擎——xNN-OCR。



揭秘支付宝中的深度学习引擎：xNN

原创：弘川 阿里技术 2017-09-28

阿里妹导读：本文介绍支付宝App中的深度学习引擎——xNN。xNN通过模型和计算框架两个方面的优化，解决了深度学习在移动端落地的一系列问题。xNN的模型压缩工具(xqeeze)在业务模型上实现了近50倍的压缩比，使得在包预算极为有限的移动App中大规模部署深度学习算法成为可能。xNN的计算性能经过算法和指令两个层面的深度优化，极大地降低了移动端DL的机型门槛。



近来，深度学习（DL）在图像识别、语音识别、

端侧AI的团队

业务算法

运用深度学习技术解决具体业务问题的专家，专精图像分类/分割/目标检测/OCR等算法

基础算法

深入理解机器学习原理、深度学习基础网络结构、信息论，主攻自动化模型压缩与加速工具的研发

工程架构

打通端侧框架 (Android, iOS) 与后台系统 (Java)，建设覆盖算法研发全流程的高效工程平台

算法优化

计算机体系结构与算法性能领域的专家，运用 NEON/OpenCL/汇编等技术提升端侧框架的效率



北京
2019

遨游“视”界 做你所想
Explore World, Do What You Want

Thank you



出品: LiveVideoStack CSDN
—— 音视频技术社区 ——