



第九届中国系统架构师大会
SYSTEM ARCHITECT CONFERENCE CHINA 2017

云端图像技术的深度学习模型与应用

李东亮

360 人工智能研究院

lidongliang@360.cn

2017.10.20



中国最大的互联网安全公司

360电脑安全产品

月活跃数达到4.42亿

360手机安全产品

移动端用户总数已达约1.49亿

360浏览器

月活跃用户数量为3.03亿

360导航

日均独立访问用户为8900万人
日均点击量约为4.51亿次

360搜索

稳定拥有35%以上的市场份额

360智能硬件

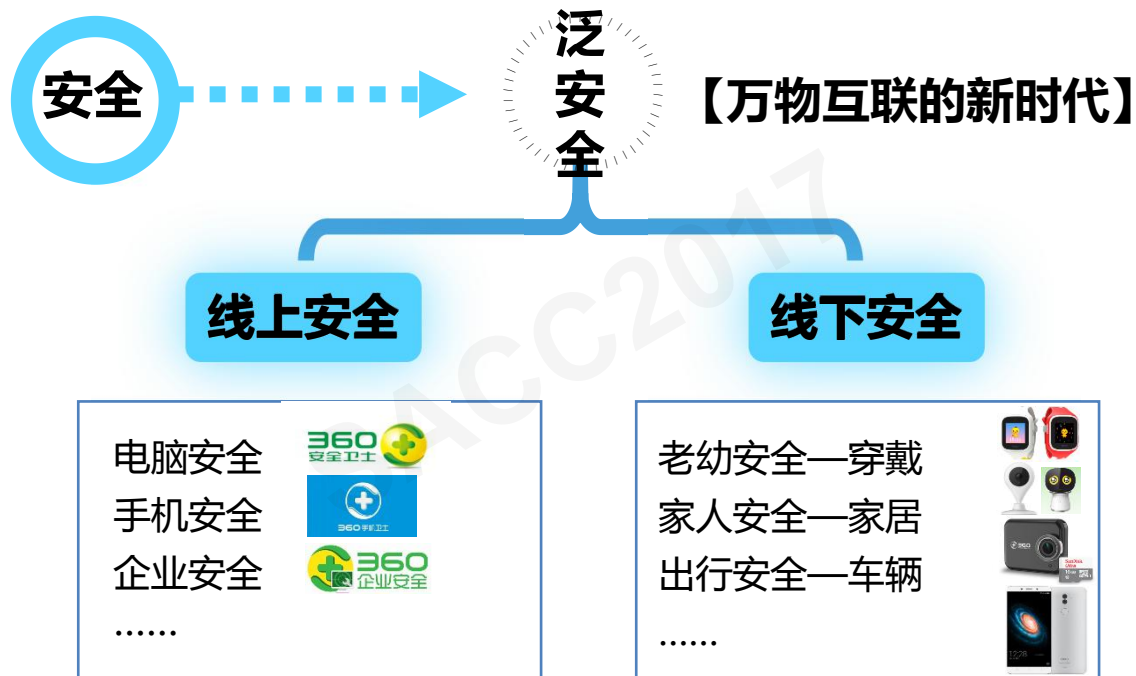
智能摄像头超400万，儿童手表超
350万，行车记录仪超300万

奇虎360



安全——360的基因

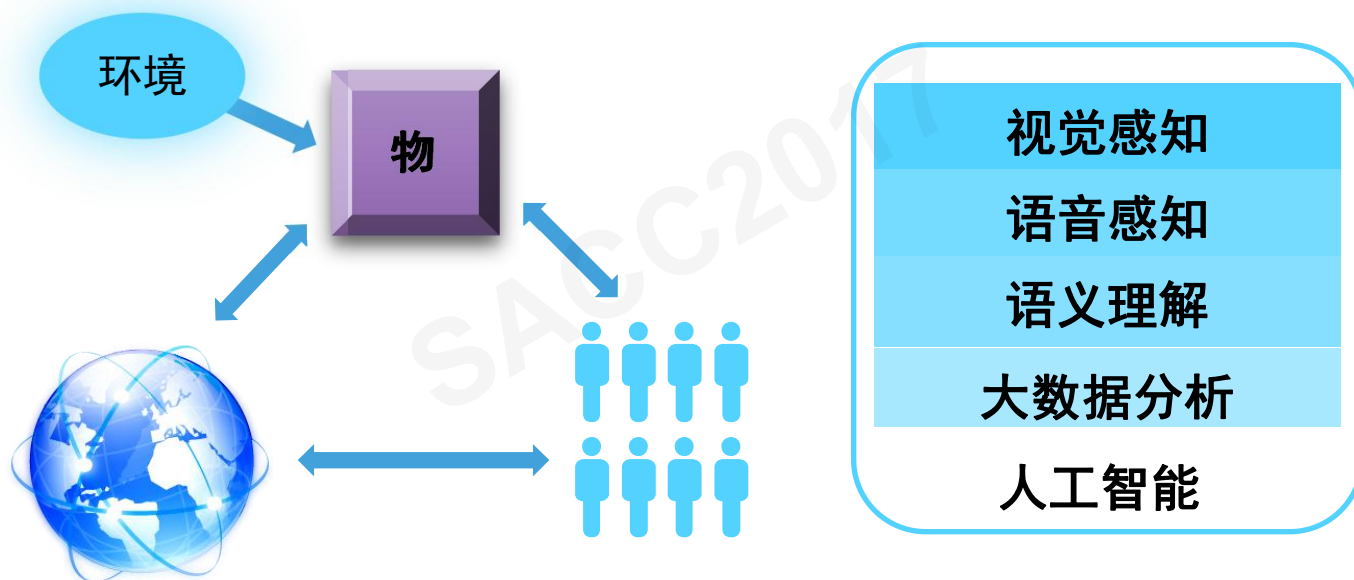
新时代的奇虎360



万物互联的新时代



万物互联的核心技术



视觉感知模型



业务

云端

移动端

数据

图像

视频

核心

检测

识别

分割

跟踪

视觉感知核心问题

核心

检测

识别

分割

跟踪

Object
Classification



Person, Horse,
Barrier, Table, etc

Object
Detection



Object
Segmentation



图像技术的三个核心难点>>小、快、准

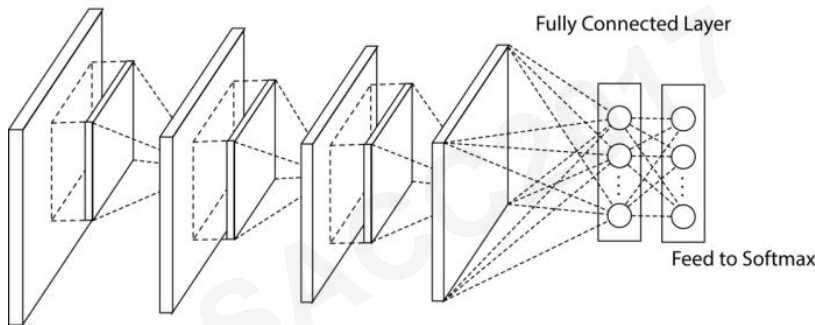
小模型



Frequent remote upgrade



Convolutional Layer



Fully Connected Layer

Feed to Softmax

线上速度快



CPU-constrained, real-time



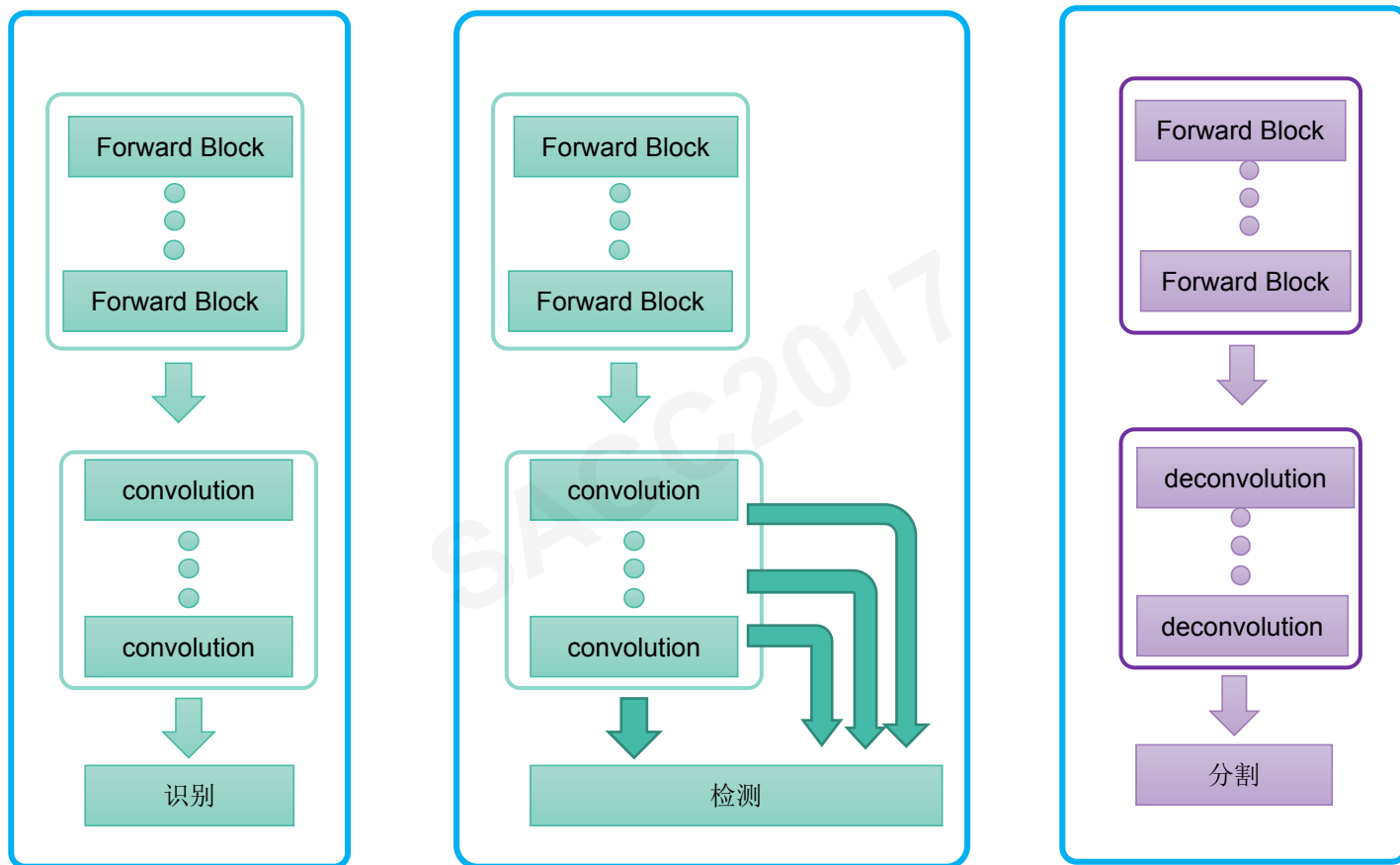
Cloud processing



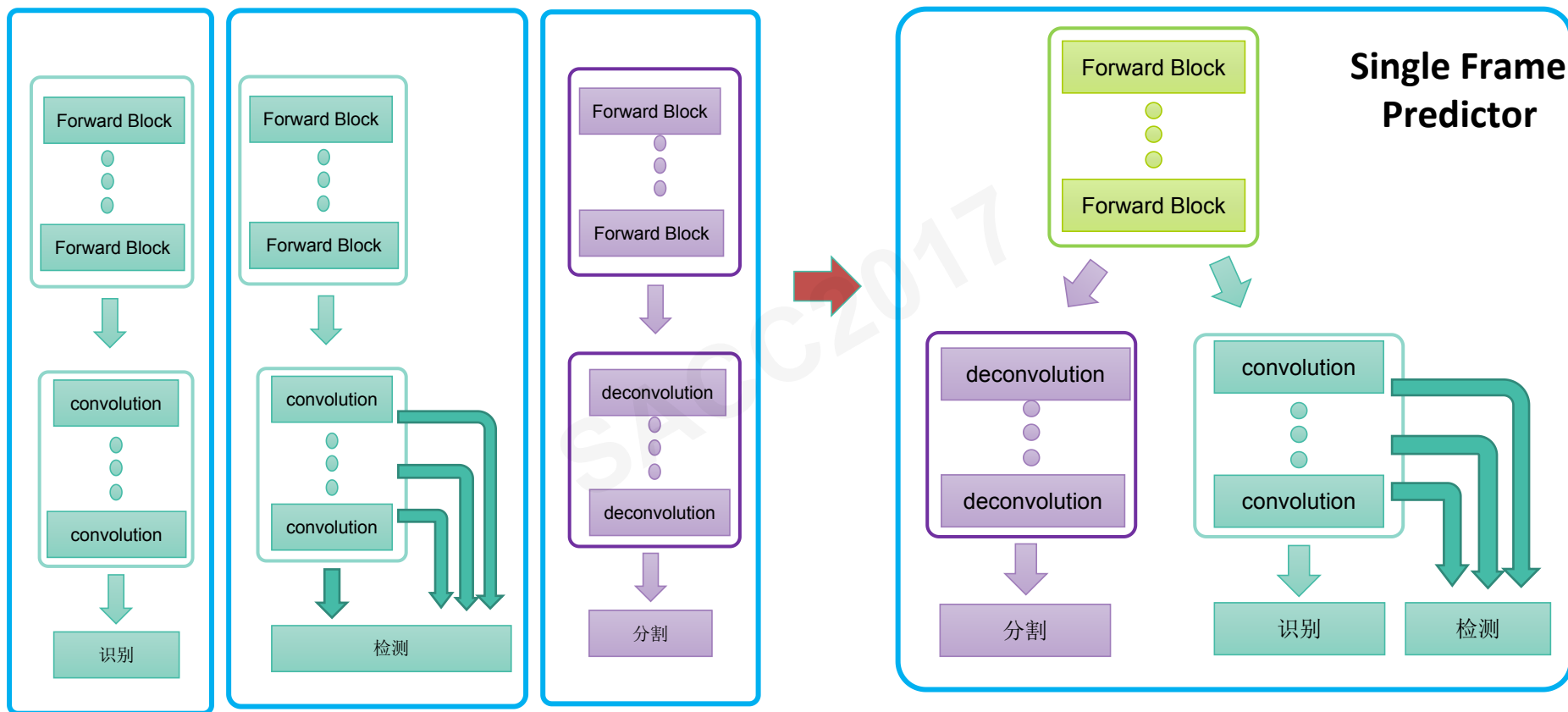
预测准

深度学习已经逐步取代各领域的传统方法

视觉感知模型



视觉感知模型-融合



视觉感知模型-融合

核心

深度学习

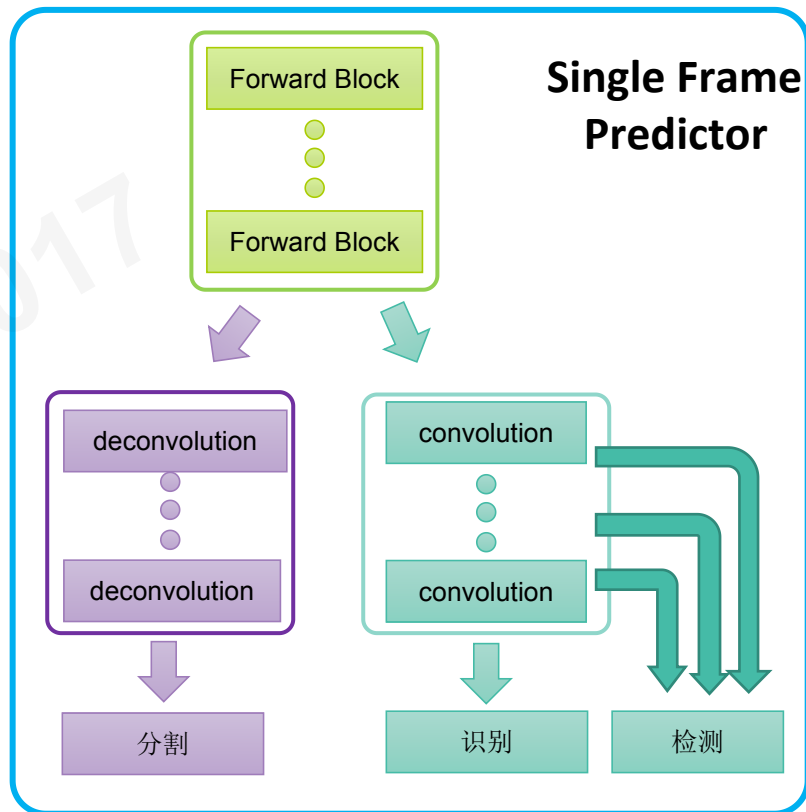
检测

识别

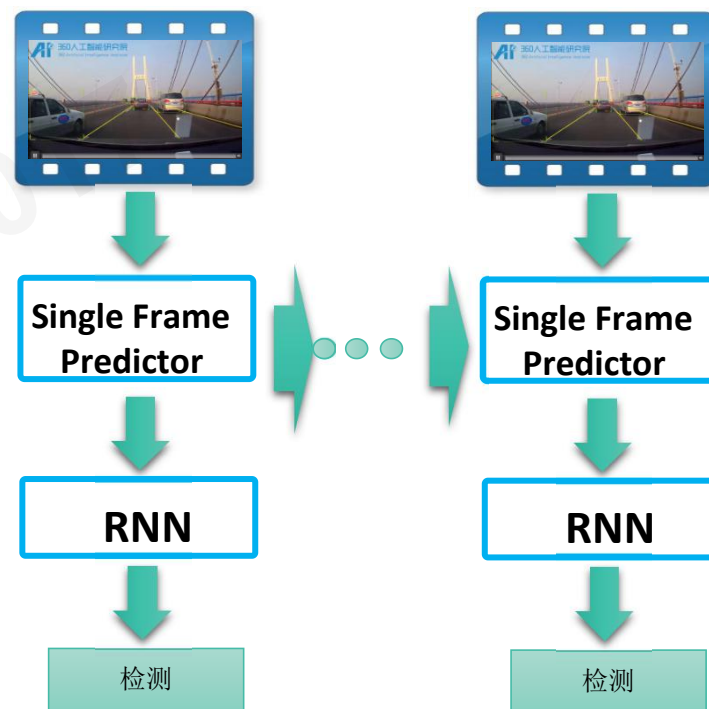
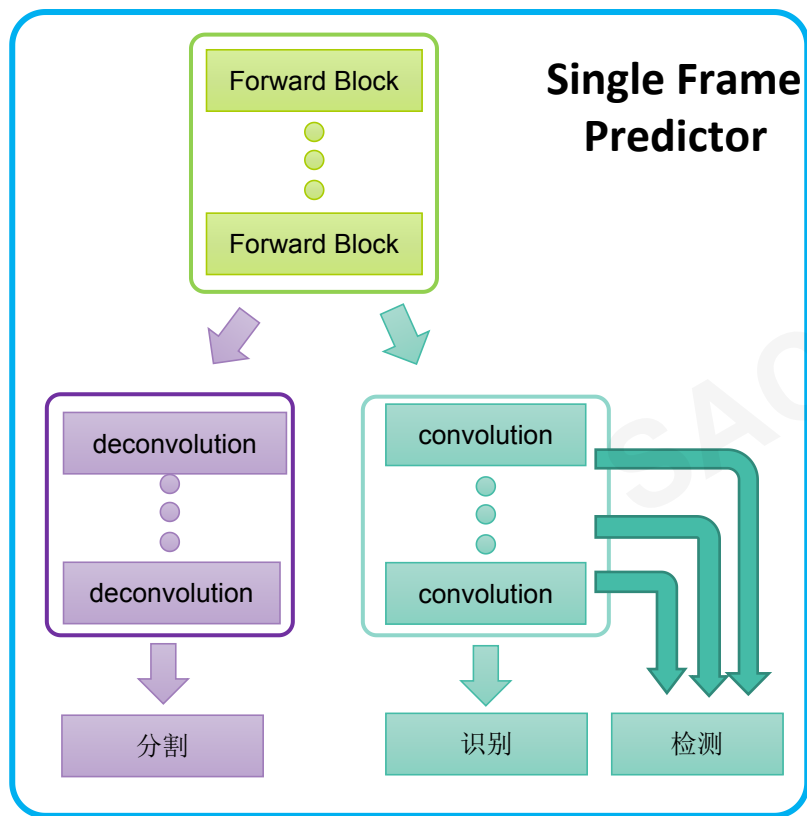
分割

跟踪

- 完全基于深度学习
- 统一分类，检测，分割，跟踪
 - ✓ 通过共享计算提高算法效率
 - ✓ 通过多个相关任务共同学习提高算法性能
- 稀疏标注
 - ✓ 在节省标注工作量的同时，充分利用视频数据



视觉感知模型-视频

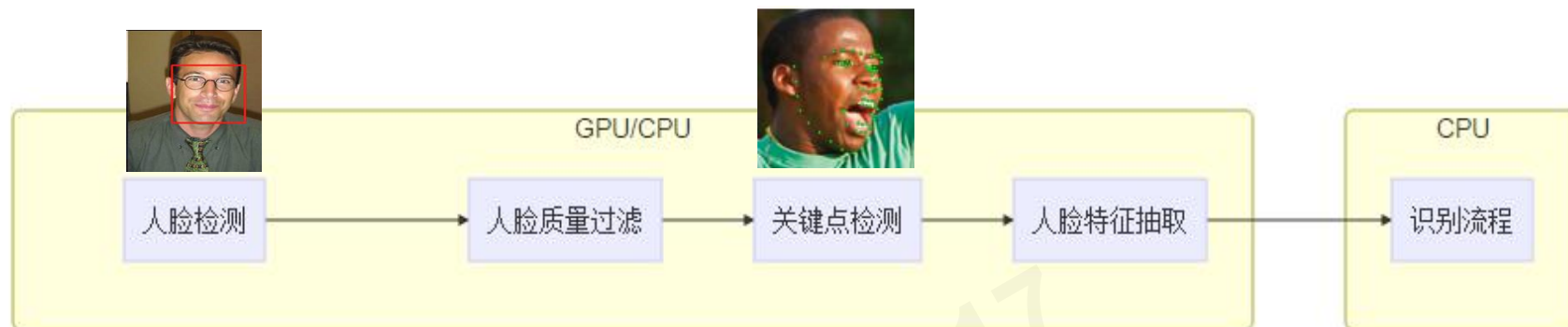


360小水滴摄像机：视觉大不同



Cloud-API 每天调用1.5亿次!2000QPS!

系统框架



- 根据业务需求，对图像人脸进行识别，将结果推送到业务端
- 基于深度学习的准确的人脸检测、特征抽取
- 人脸检测占用95%计算资源
- 峰值时会达到1500 QPS

检测-人脸检测/人形检测

场景多样、人脸小、位置边缘



本页图片均来自公开摄像头

检测-人脸检测/人形检测



	手机	服务器
可缩小尺寸	240P	720P
CPU	ARM（千元机）	E5-2630
时间	50ms	120ms
GPU		2-5ms(K40)

图像技术的三个核心难点>>小、快、准

小模型

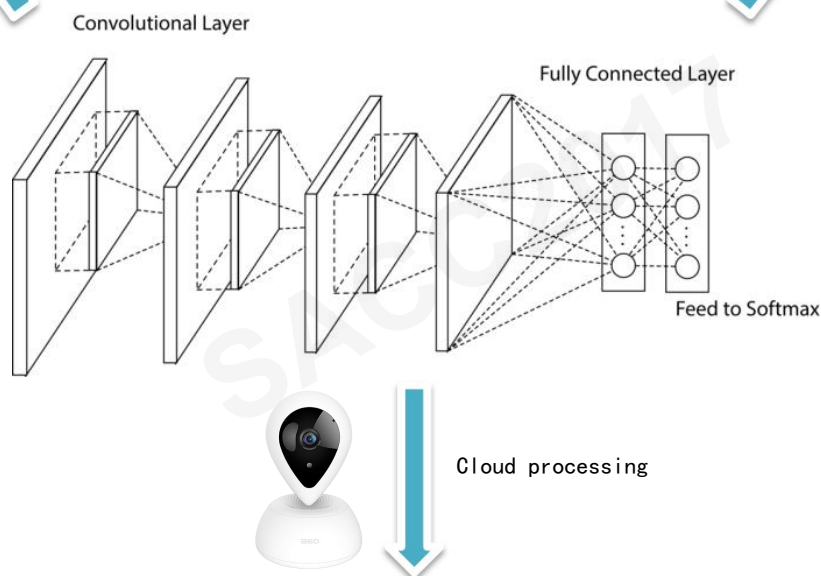


Frequent remote upgrade

线上速度快



CPU-constrained, real-time



图像技术的三个核心难点>>小、快、准

数据

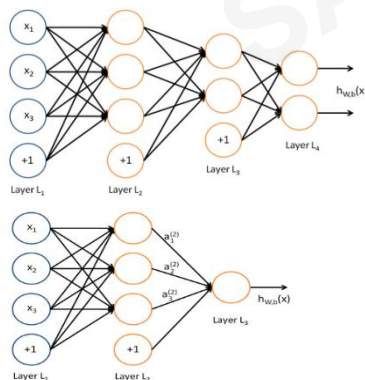


工程

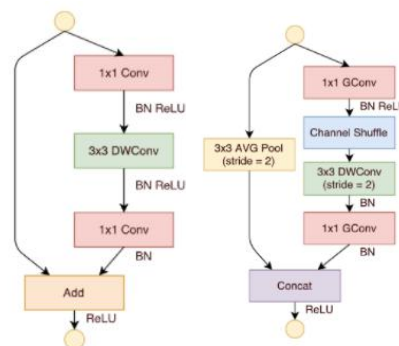


模型

模型缩减

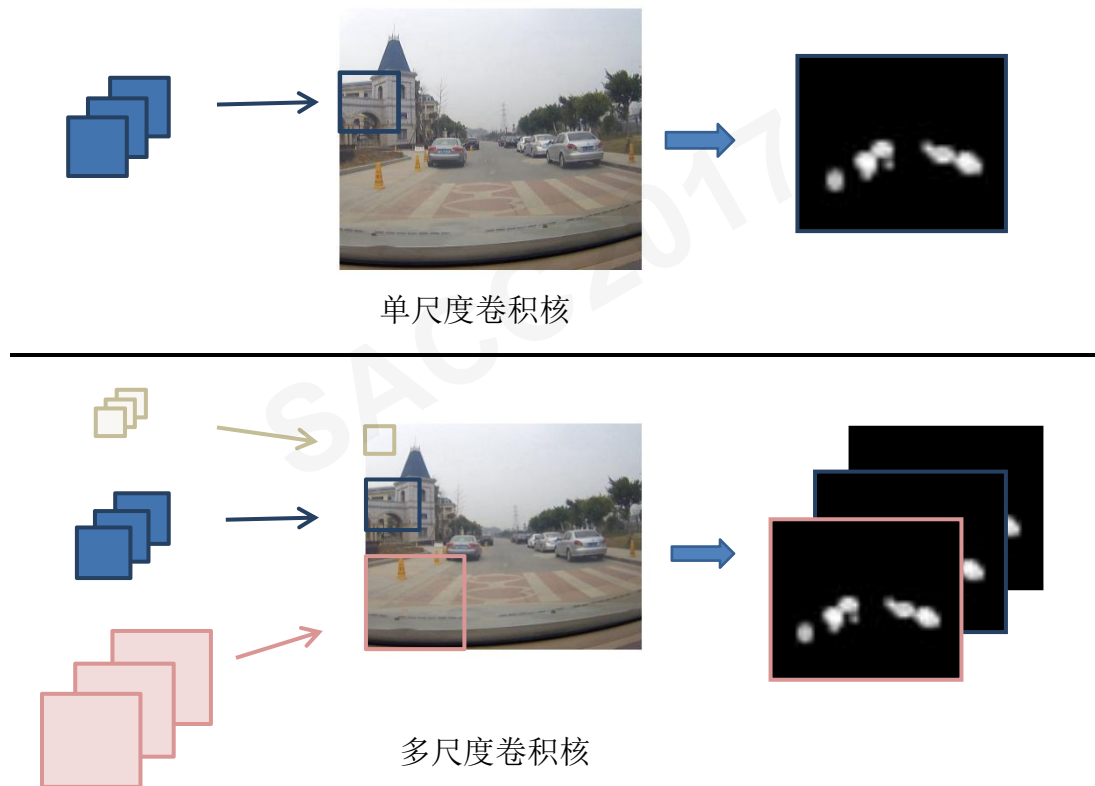


结构演进



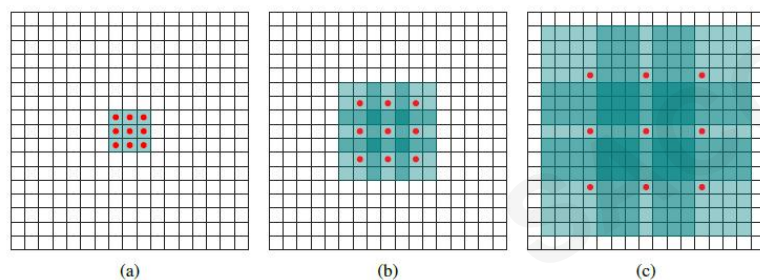
视觉感知的三个核心难点>>小、快、准

Inception结构

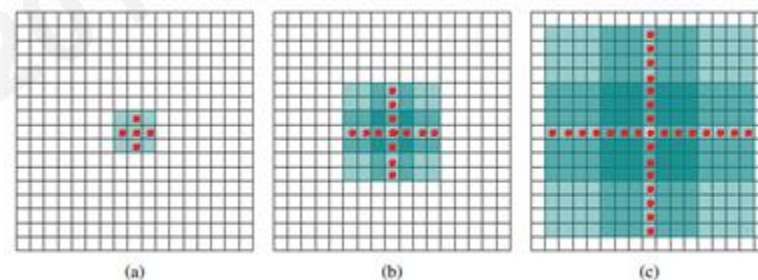


视觉感知的三个核心难点>>小、快、准

稀疏卷积核



Hole algorithm



Cross-convolution

视觉感知的三个核心难点>>小、快、准

低秩矩阵分解

- 复杂度分析

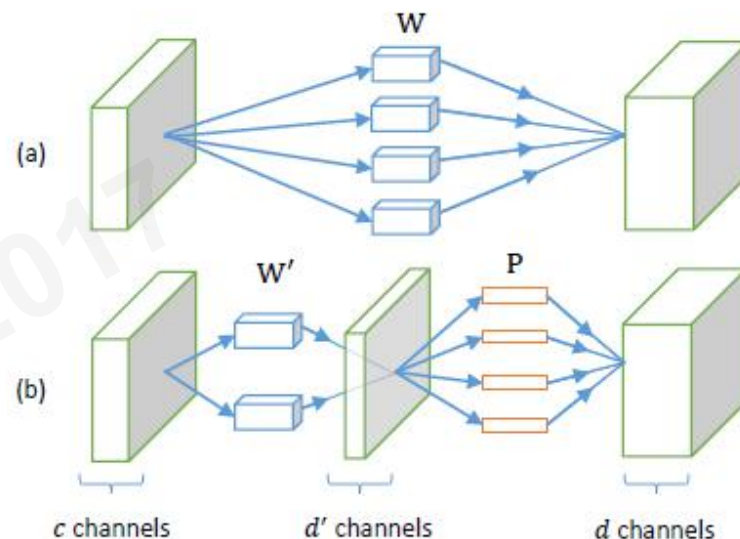
$$y = Wx$$

$$O(dk^2c)$$

$$y = PW'x + b$$

$$O(d'k^2c) + O(dd')$$

- 问题求解



$$\min_M \sum_i \|(y_i - \bar{y}) - M(y_i - \bar{y})\|_2^2,$$
$$s.t. rank(M) \leq d'$$



**家人/陌生人识别
疲劳监控
萌拍、换脸等娱乐/游戏功能**

**准确、稳定、鲁棒、低功耗
人脸检测、定位、识别多项世界前沿的性能**



最早在人脸标准库上LFW达到99.7%的团队之一!

GPU服务框架-图像特点

通用计算 (Caffe/Tensorflow/Mxnet)

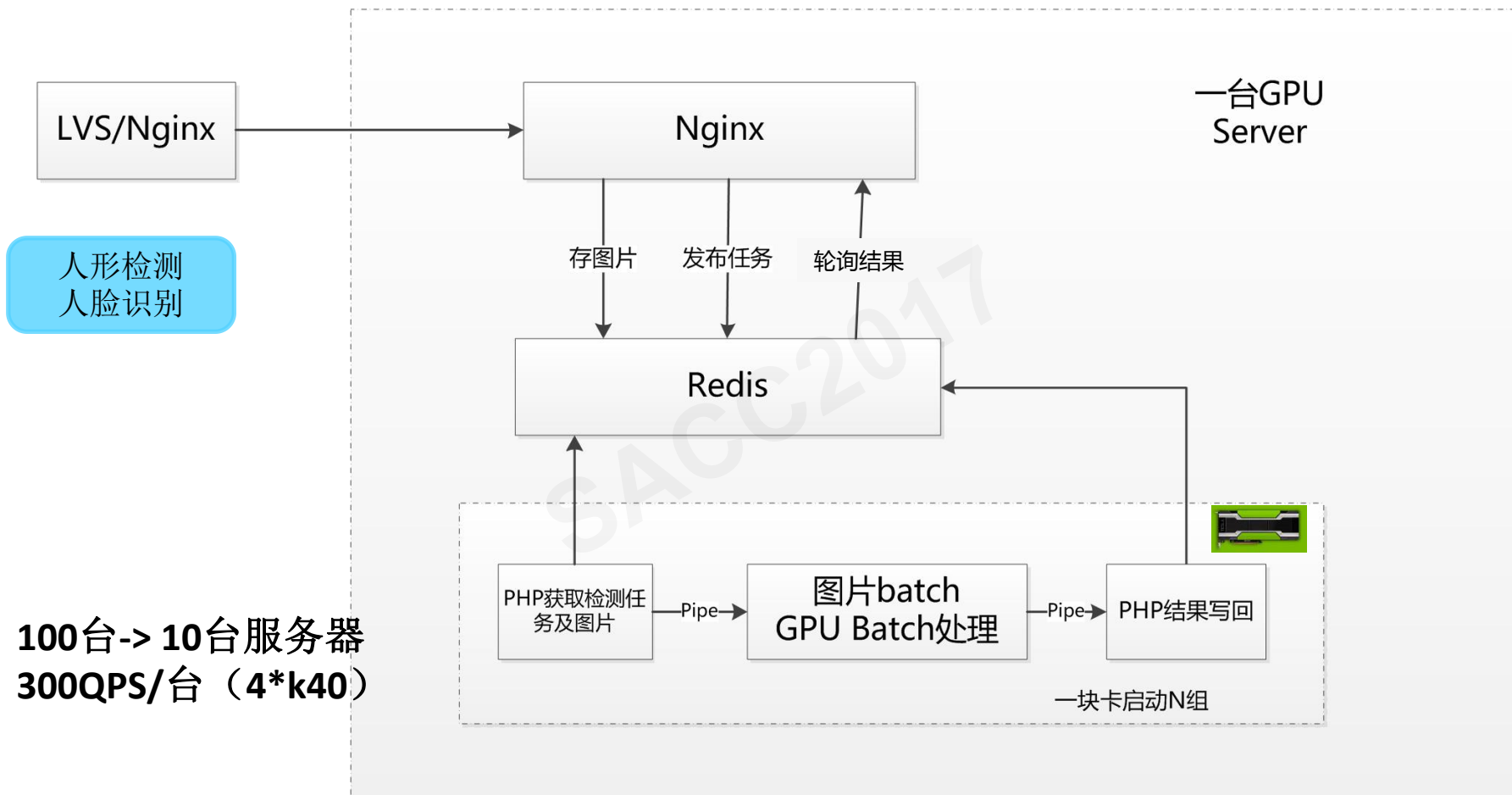
输入输出固定, 无状态

计算量大、响应→GPU

多任务串联

传输、存储压力

GPU服务框架



THANKS

