

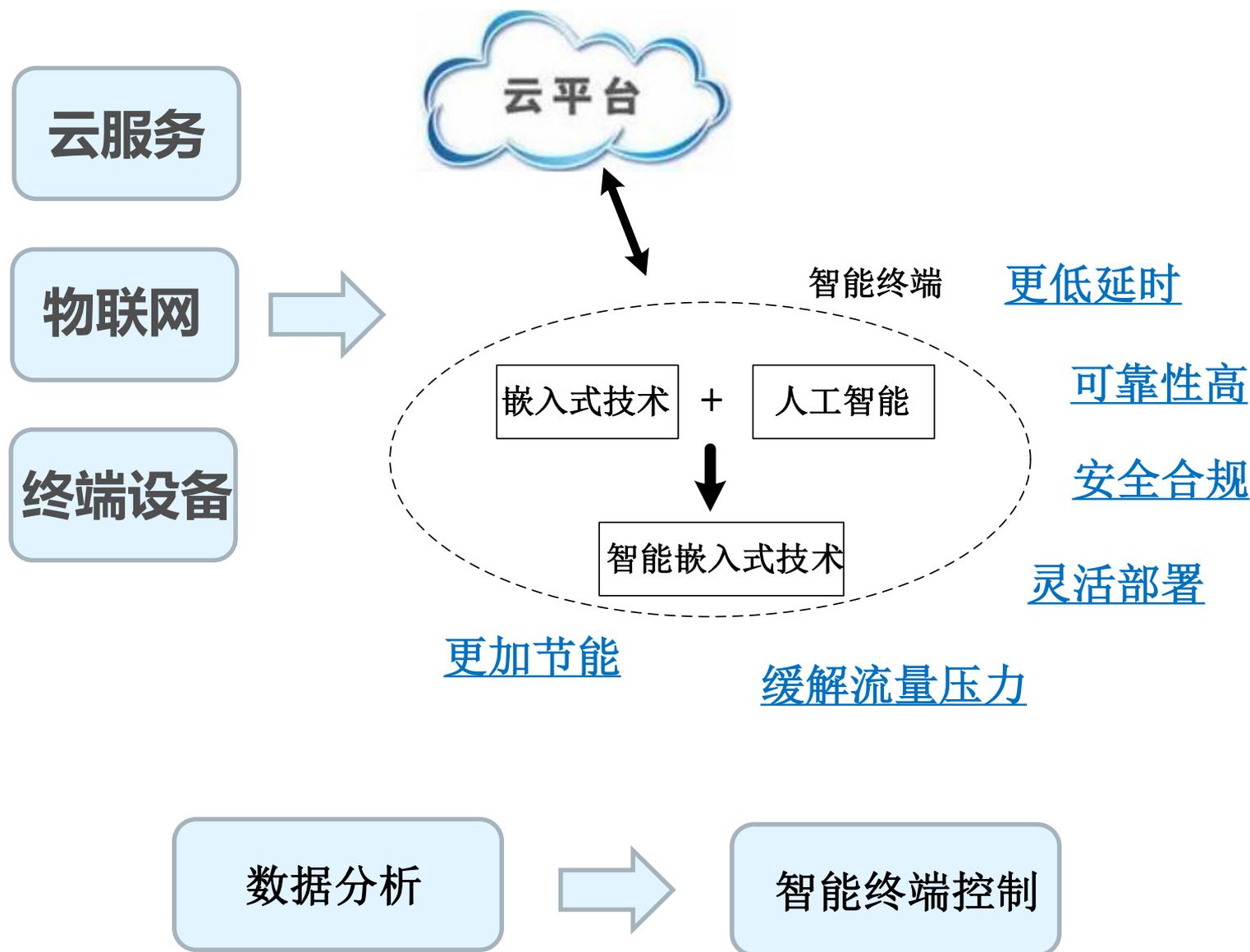


2019年中国嵌入式技术大会
EMBEDDED TECHNOLOGY
Conference China 2019

智能嵌入式技术及应用开发

华南理工大学
计算机科学与工程学院
毕盛
2019年12月19日

嵌入式系统 + 人工智能



智能嵌入式应用



结合语义的目标识别

智能抓取



自主导航定位追踪系统



自主行走



深度学习

视觉追踪

运动控制

全向运动



TurboX 智能套件

多机协作

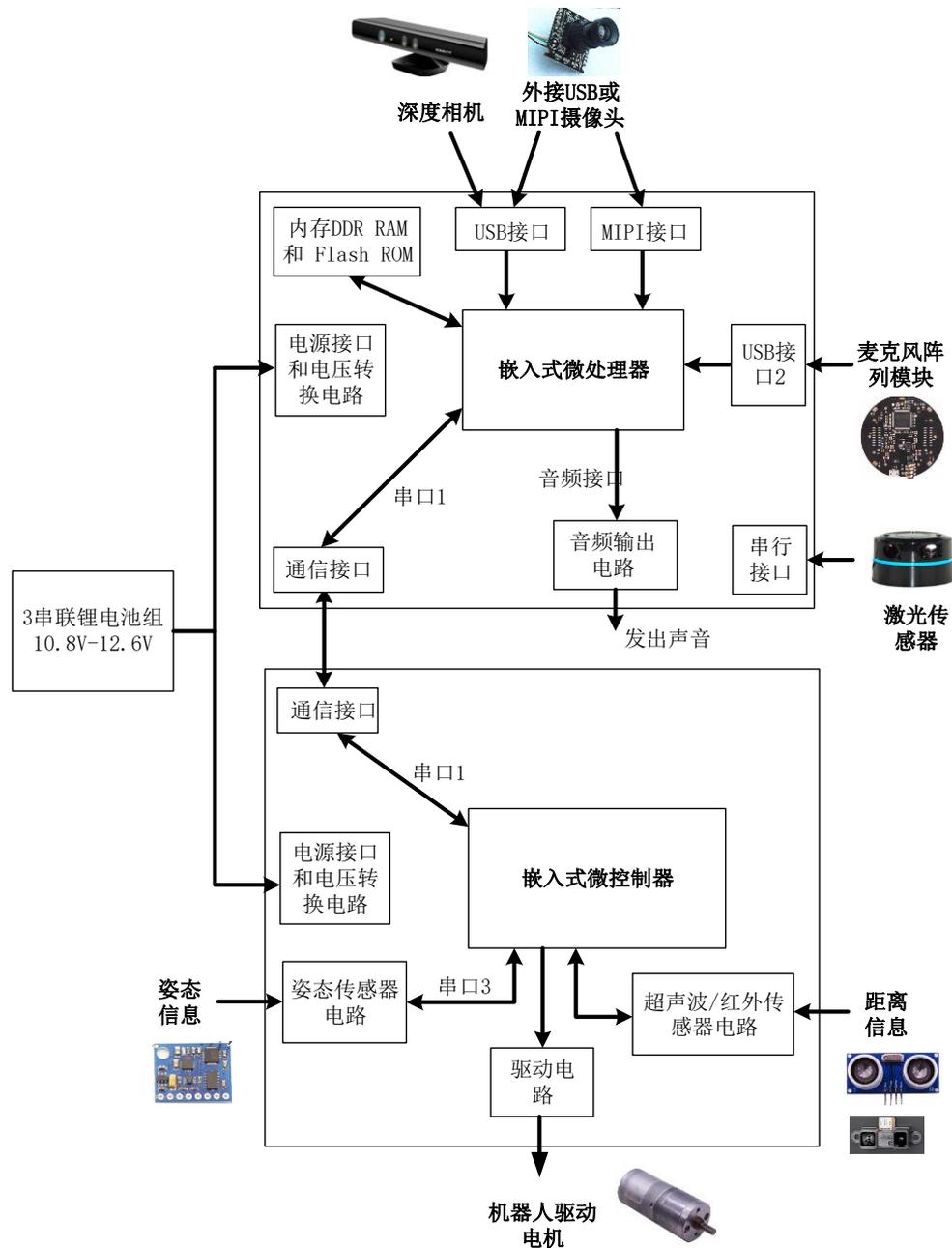
全向移动机械臂



SLAM技术

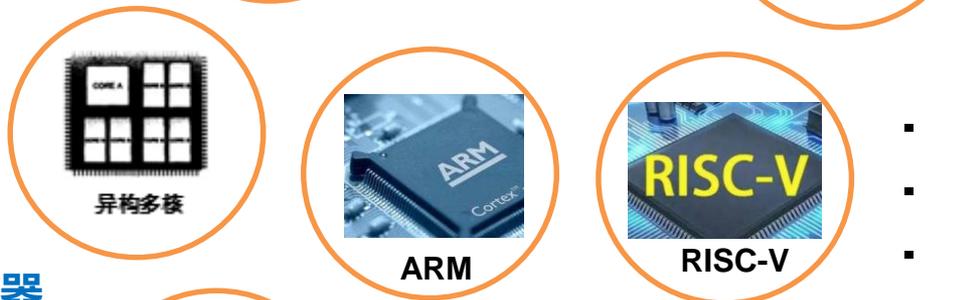
智能控制及强化学习

智能嵌入式通用架构

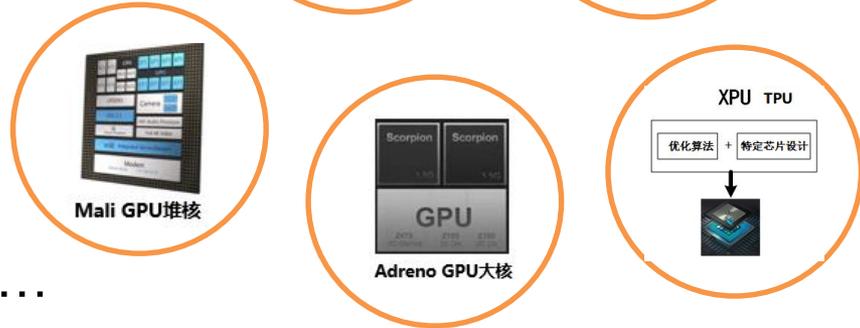


智能嵌入式硬件芯片

MPU-嵌入式微处理器



MCU-嵌入式微控制器



DSP-数字信号处理器



SOC-片上系统



(神经网络芯片) NN



Open AI LAB EAIDK平台



中科创达 TurboX AI Kit



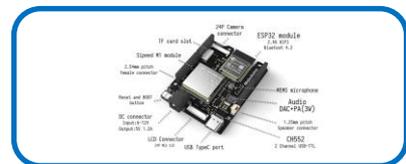
NVIDIA Jetson 平台



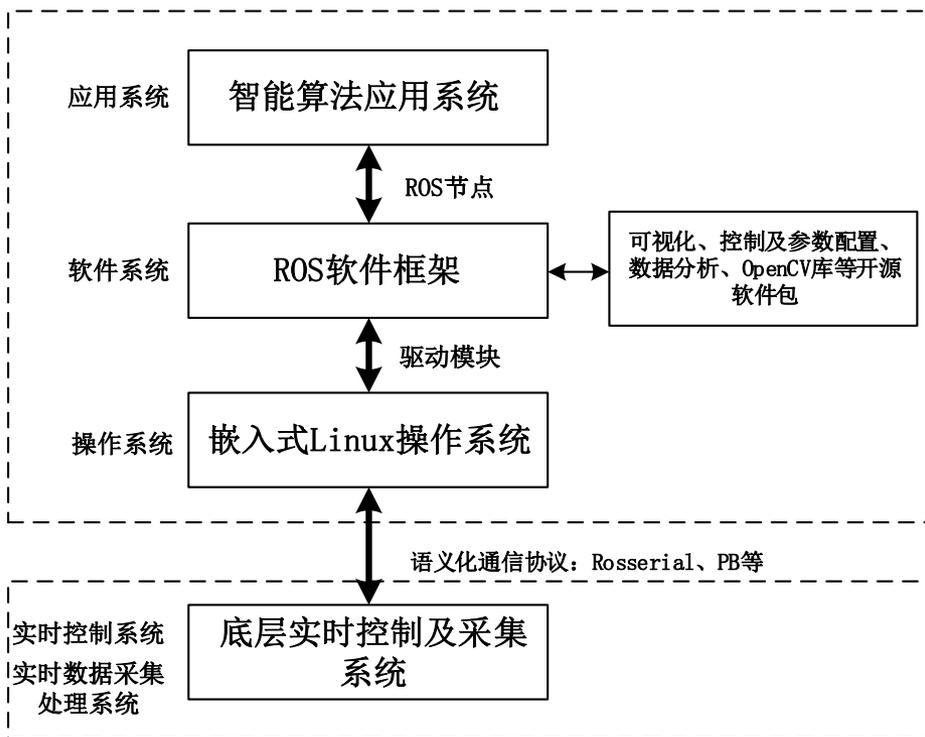
华为Atlas系列



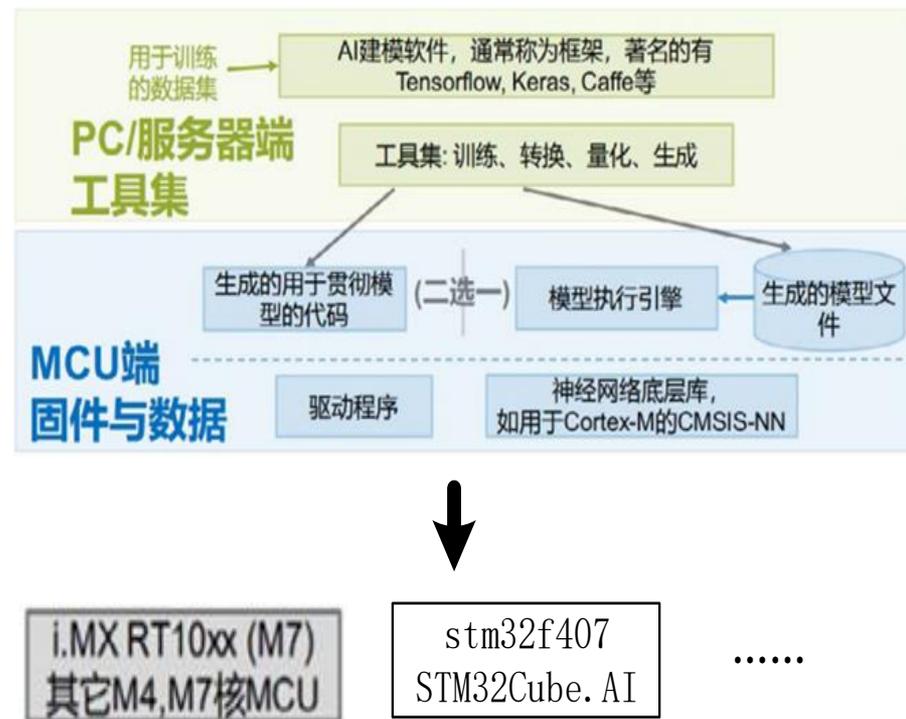
Kendryte K210



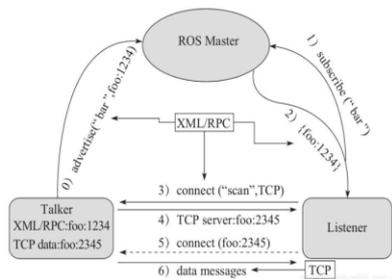
MPU智能嵌入式软件系统



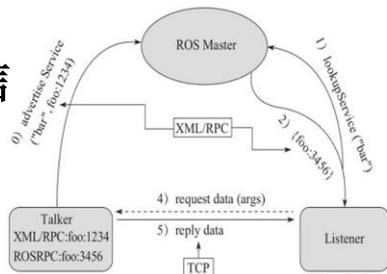
MCU智能嵌入式软件系统



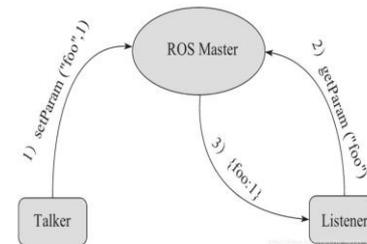
话题通信机制



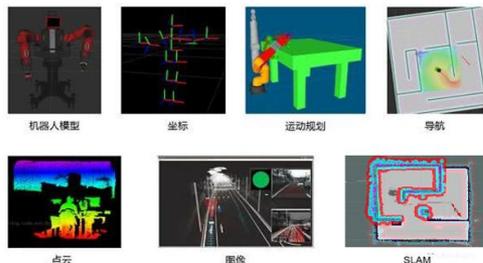
服务通信机制



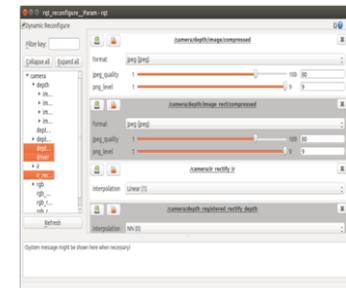
共享参数机制



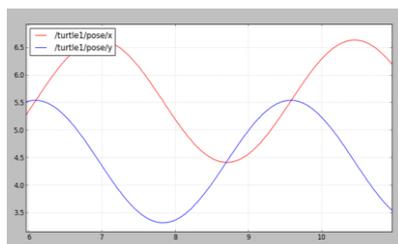
rviz可视化工具



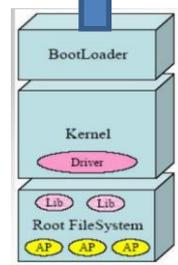
参数动态配置



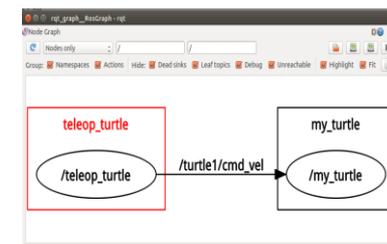
数据绘图

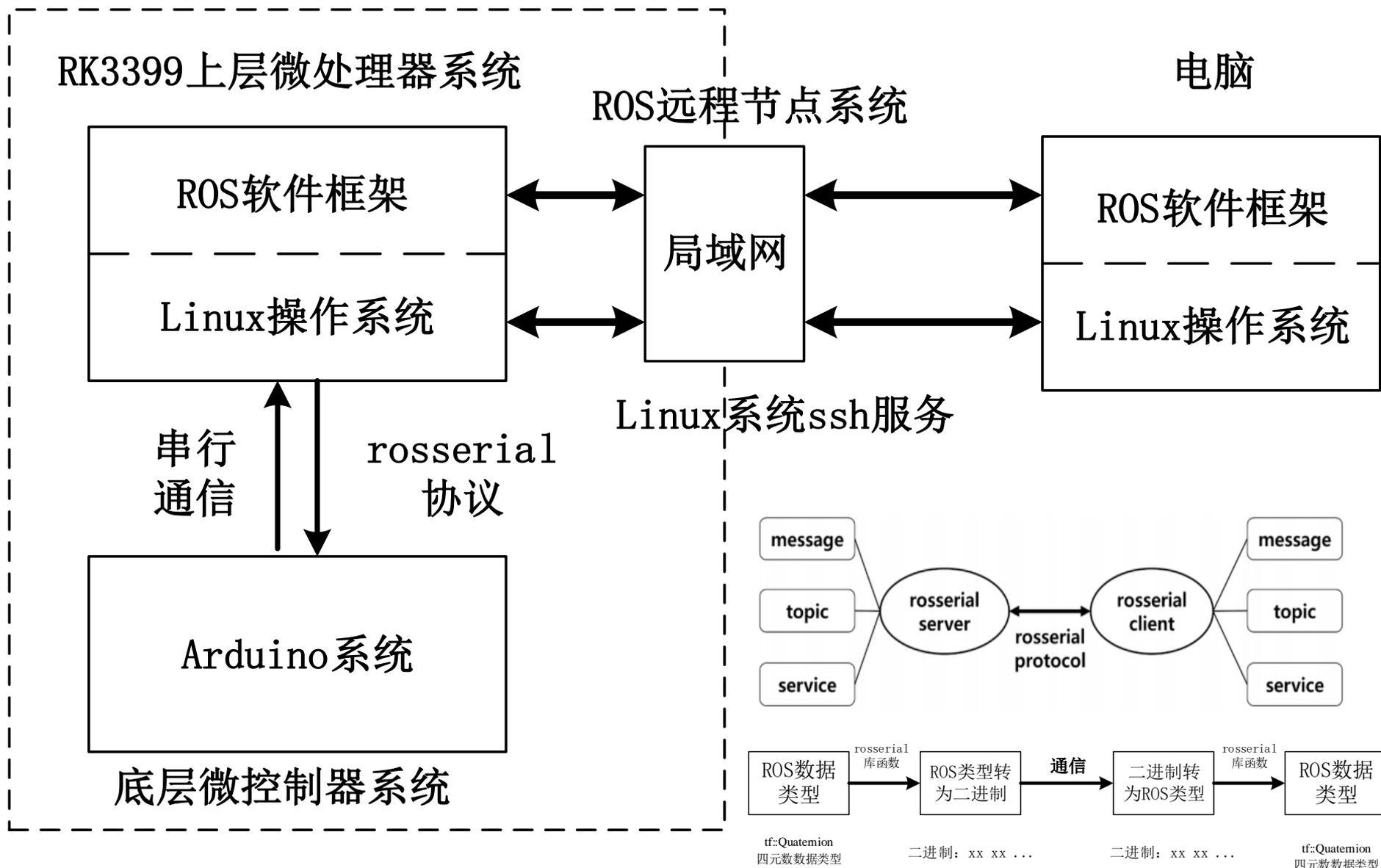


操作系统

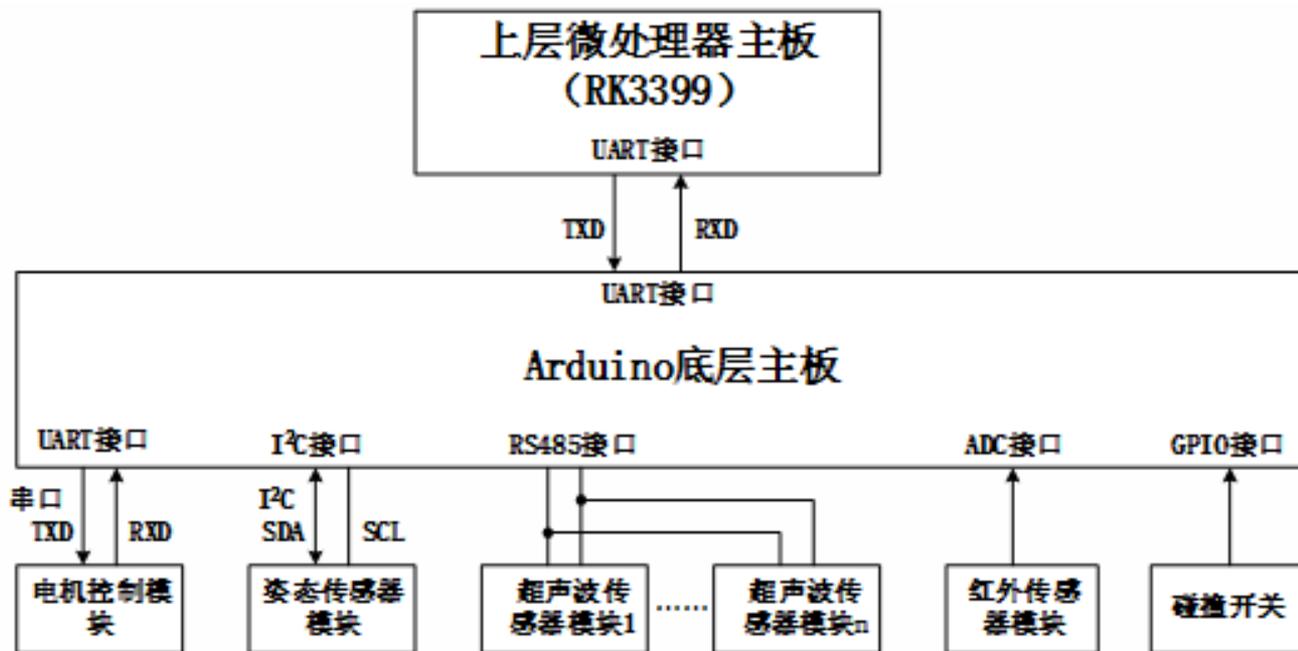
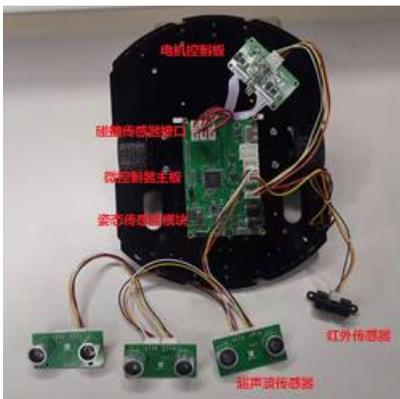


计算图可视化

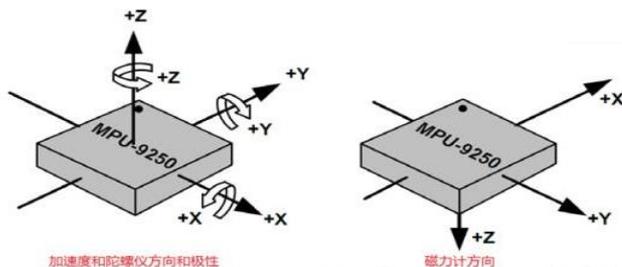




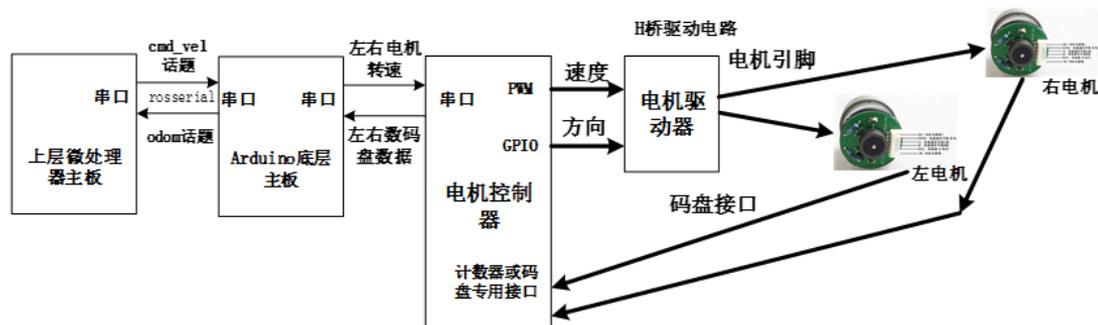
底层智能传感器系统

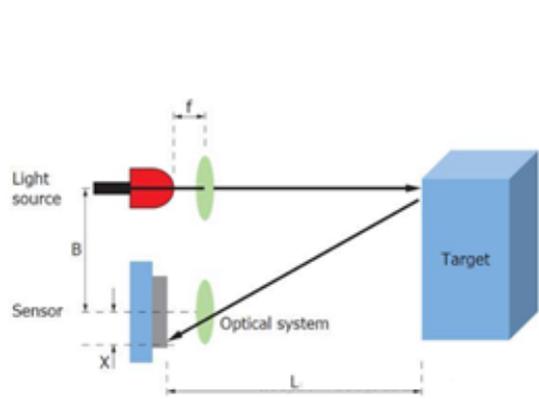


9轴融合:陀螺仪+加速度+电子罗盘四元数

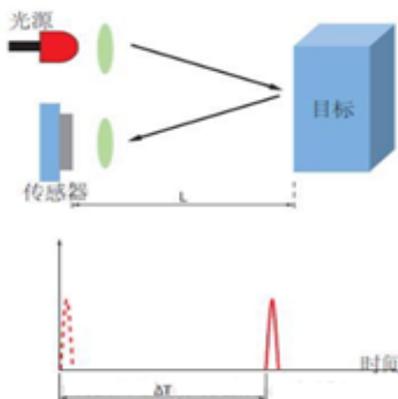


机器人小车控制

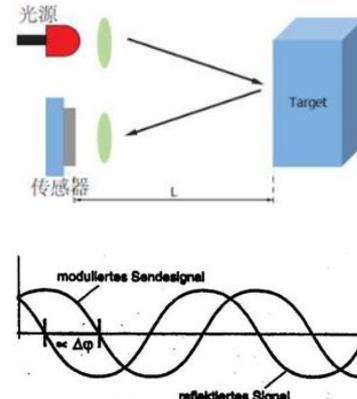




光学三角激光测距

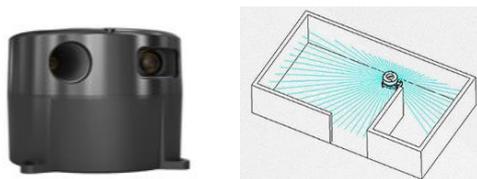


脉冲式TOF

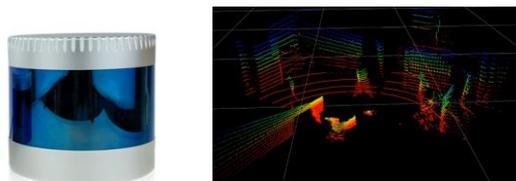


相位式TOF

单线激光



多线激光



TOF传感器



扫描匹配方法

卡尔曼滤波

迭代最近点
(ICP)

Hector方法

概率栅格方法

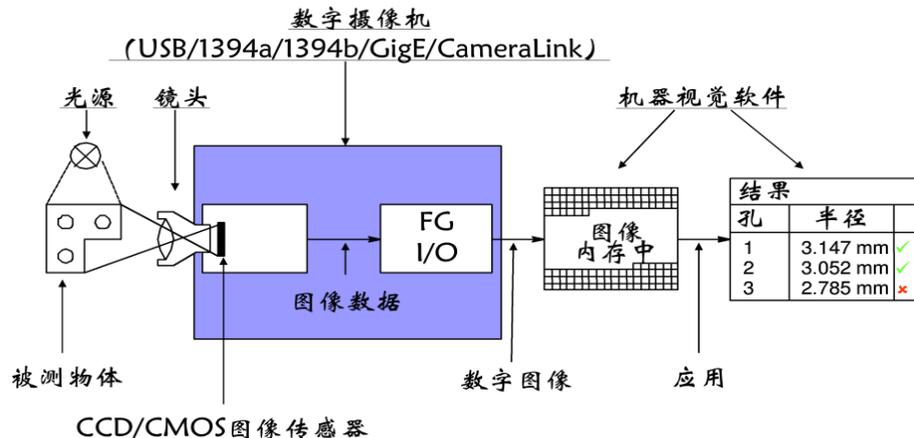
粒子滤波

正态分布变换
(NDT)

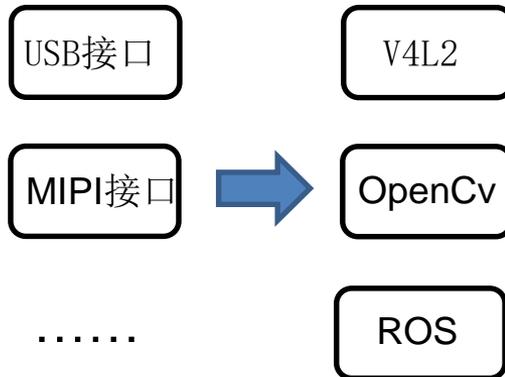
Cartographer
方法

视觉智能传感器系统

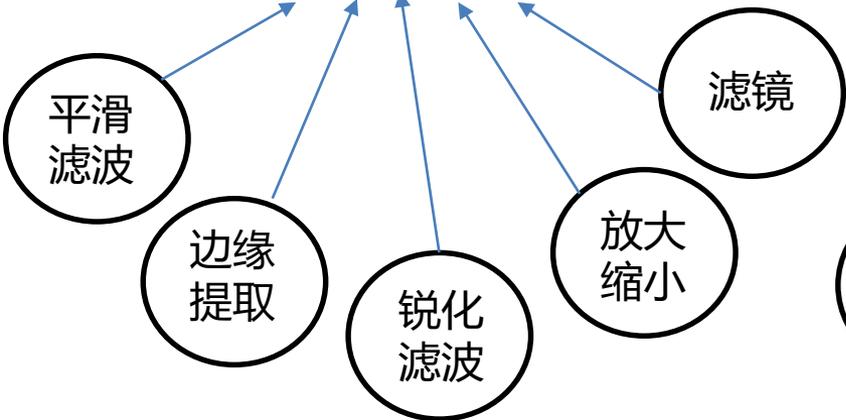
机器视觉系统构成——数字摄像机



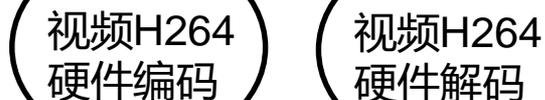
摄像头



图像处理



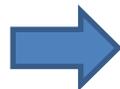
视频处理



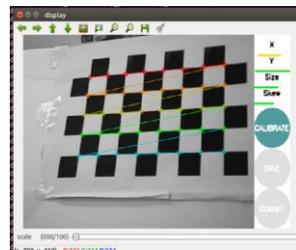
机器视觉



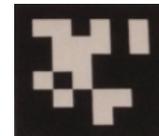
视觉定位系统



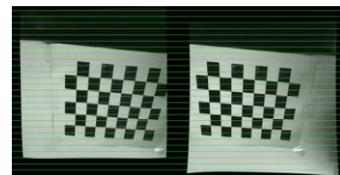
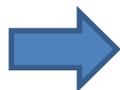
标定



标签定位



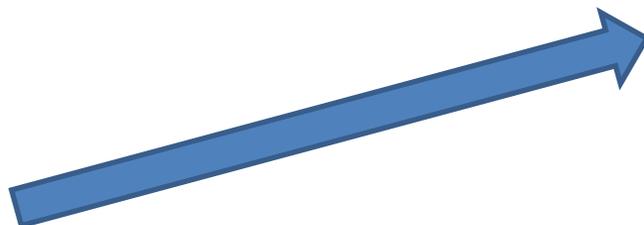
双目相机



距离测量



RGBD深度相机

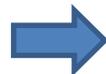


视觉SLAM

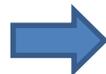
ORB和LSD方法等

自然语言系统

麦克风



语音采集



前处理模块

音频编解码

噪声抑制

语音活性检测

回声消除



语音识别

特征提取器

语音唤醒技术

语音识别

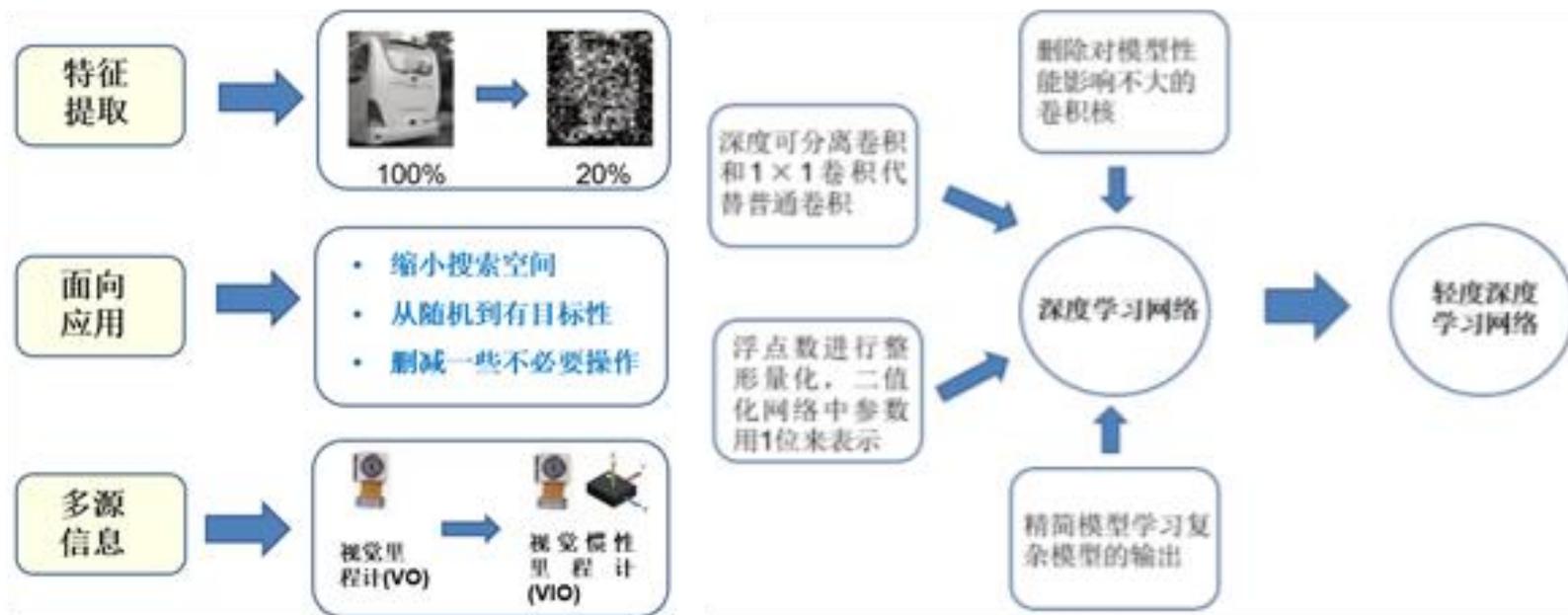
语音合成



自然语言理解



深度学习部署



在深度学习过程中，需要对网络进行简化，主要包括：1) 删除对模型性能影响不大的卷积核；2) 深度可分离卷积和 1×1 卷积代替普通卷积；3) 浮点数进行整形量化，二值化网络中参数用1位来表示；4) 精简模型学习复杂模型的输出。

一些公司也开发了前端部署方案，例如ARM公司OPEN AI LAB的Tengine框架、Google公司的TensorFlow Lite、腾讯公司的NCNN框架，小米公司的MACE框架、和百度公司的Mobile-deep-learning、亚马逊公司的TVM和美国高通公司的SNPE等通过借助多核和加速单元实现卷积的快速计算，从而在移动设备上有效的实现深度学习算法。

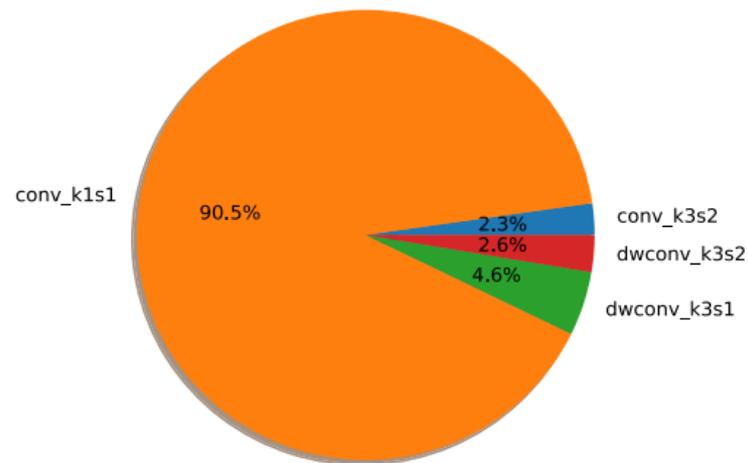
对某个特定网络的优化步骤：

1. 将一些简单的层与前一层合并，减少访存如 BatchNorm层、激活层(ReLU)。
2. 测试网络中所有网络层的推理耗时。
3. 针对耗时非常严重的网络层进行专门优化 (例如右图中的1x1卷积)。

3.1 对外层for循环，使用OpenMP将计算任务分配到多个CPU上。

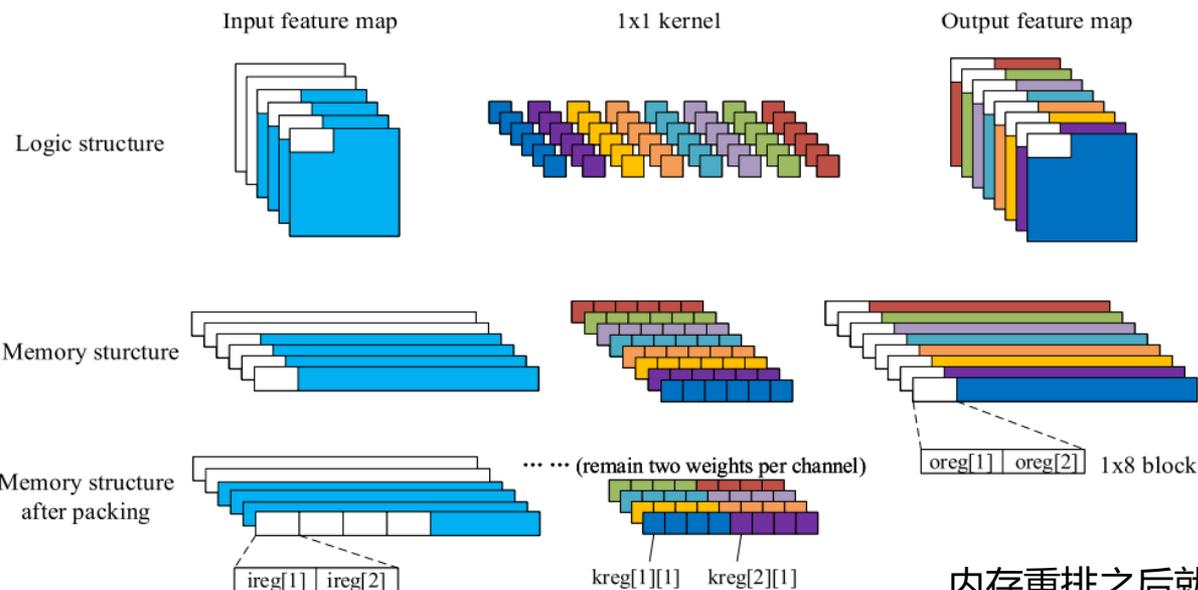
3.2 对内存for循环，先进行循环展开，然后使用ARM NEON指令来优化单个CPU上的计算。

下图是MobileNet网络的耗时测试结果，图中并没有画出softmax和池化层，因为这些层的影响很小。



1x1卷积优化举例

对于输入数据量比较大的网络层，因为数据大，那么计算过程中所需要访问的内存就会大。而对于内存的访问模式我们是可以根据计算过程来提前了解到的。所以我们可以将计算的输入数据（输入特征图和卷积核）的内存根据计算过程中的内存访问模式来重排，增强访存的局部性，提高cache命中率，减少内存延迟对CPU计算的影响。



1x1卷积的优化方法：

最上面分别是：6个输入特征图、8个1x1卷积核、8个输出特征图的逻辑结构。

中间分别是：对应的在内存中的存储结构。

最下面是：经过内存重排后的内存布局。

内存重排之后就可以利用OpenMP和ARM NEON指令来优化1x1卷积的计算过程。针对特定的网络优化过的网络层，对于其他使用到该层的网路也会有一定的加速作用。

SLAM方法 -- 创建地图、导航及追踪

地图创建方法: 例如Gmapping,Hector, Cartographer方法等。
定位中匹配方法: 例如粒子滤波, 迭代梯度下降(ICP, NDT)等。
导航方法: 全局导航A*, D*方法, 局部导航DWA, TEB方法等。

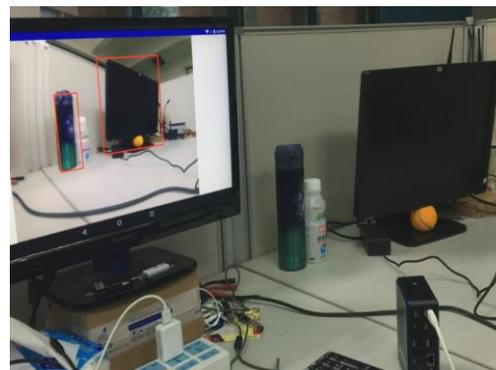


实例2-结合自然语言理解和环境认知的智能服务系统

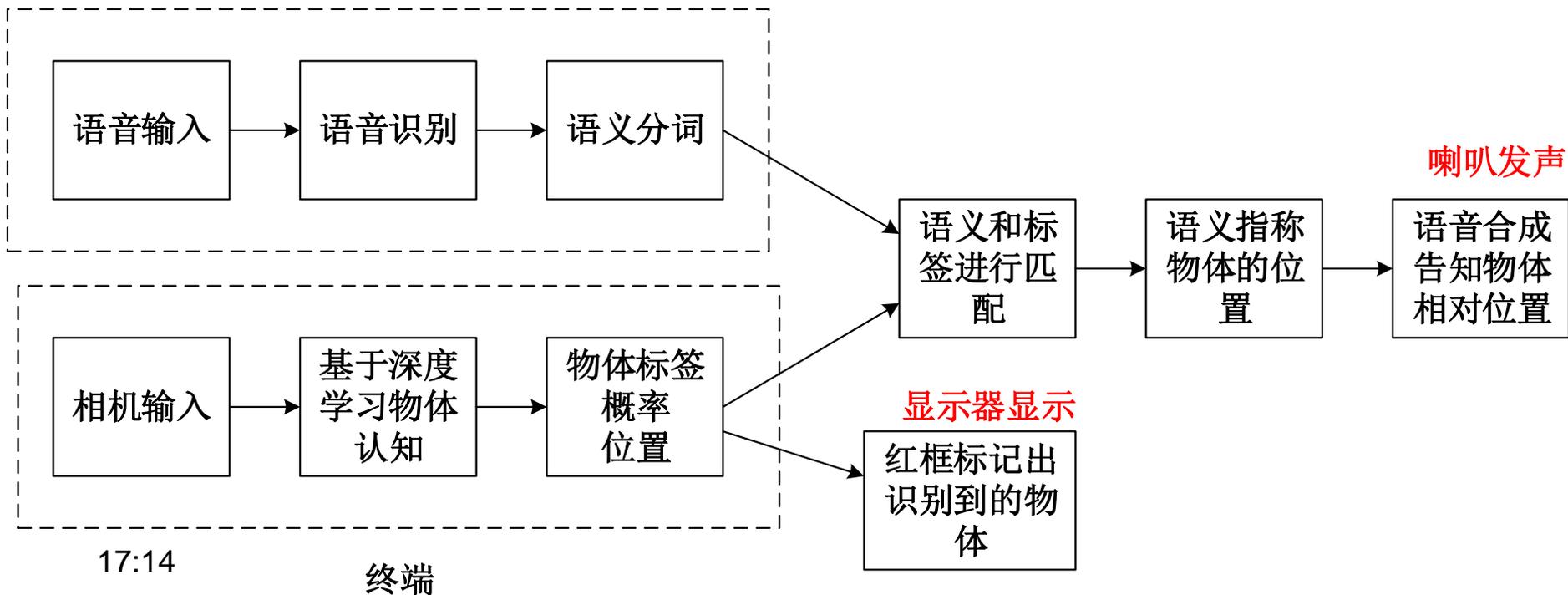
自然语言
理解方法

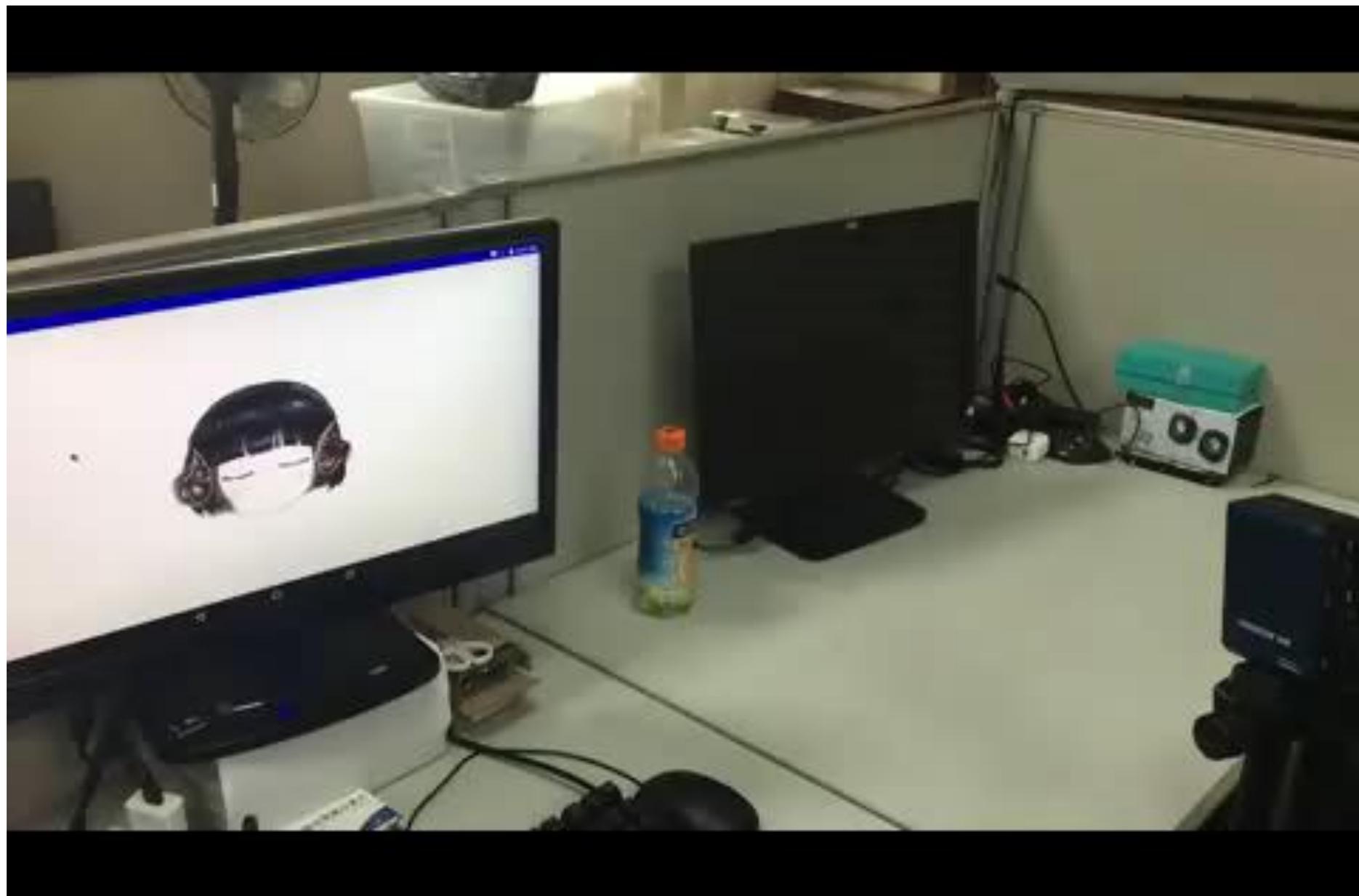


基于深度
学习物体
认知方法



云端





■ 谢谢!

