

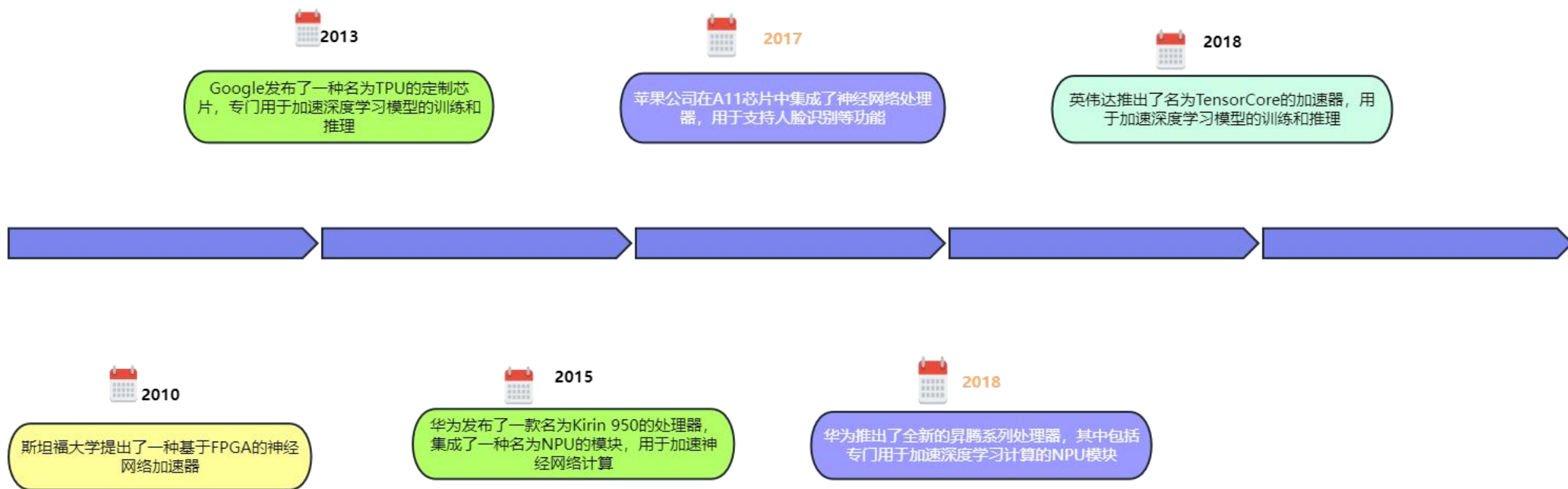
RKNPU教程

01 初识RKNPU

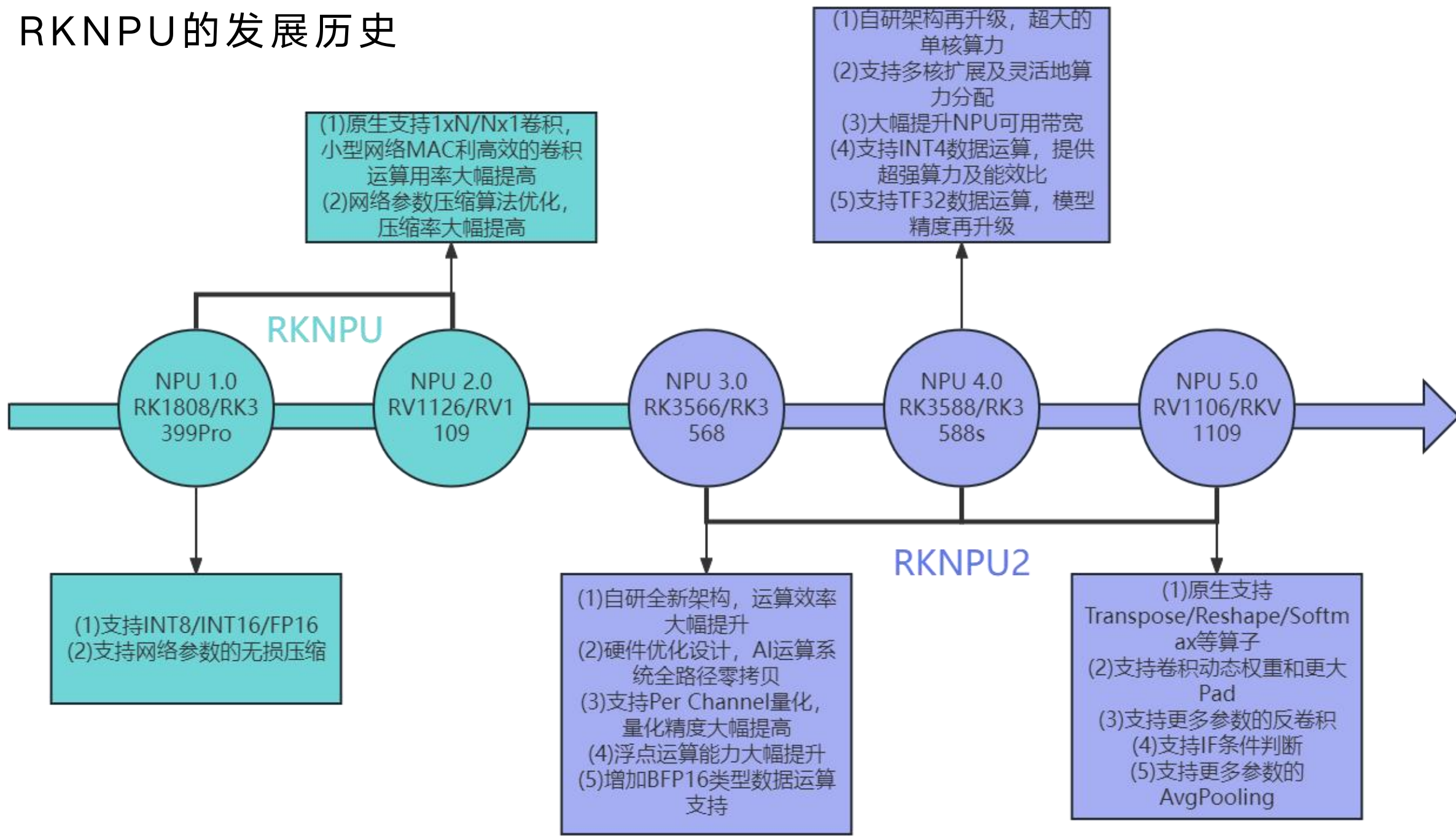
TOPEET 迅为



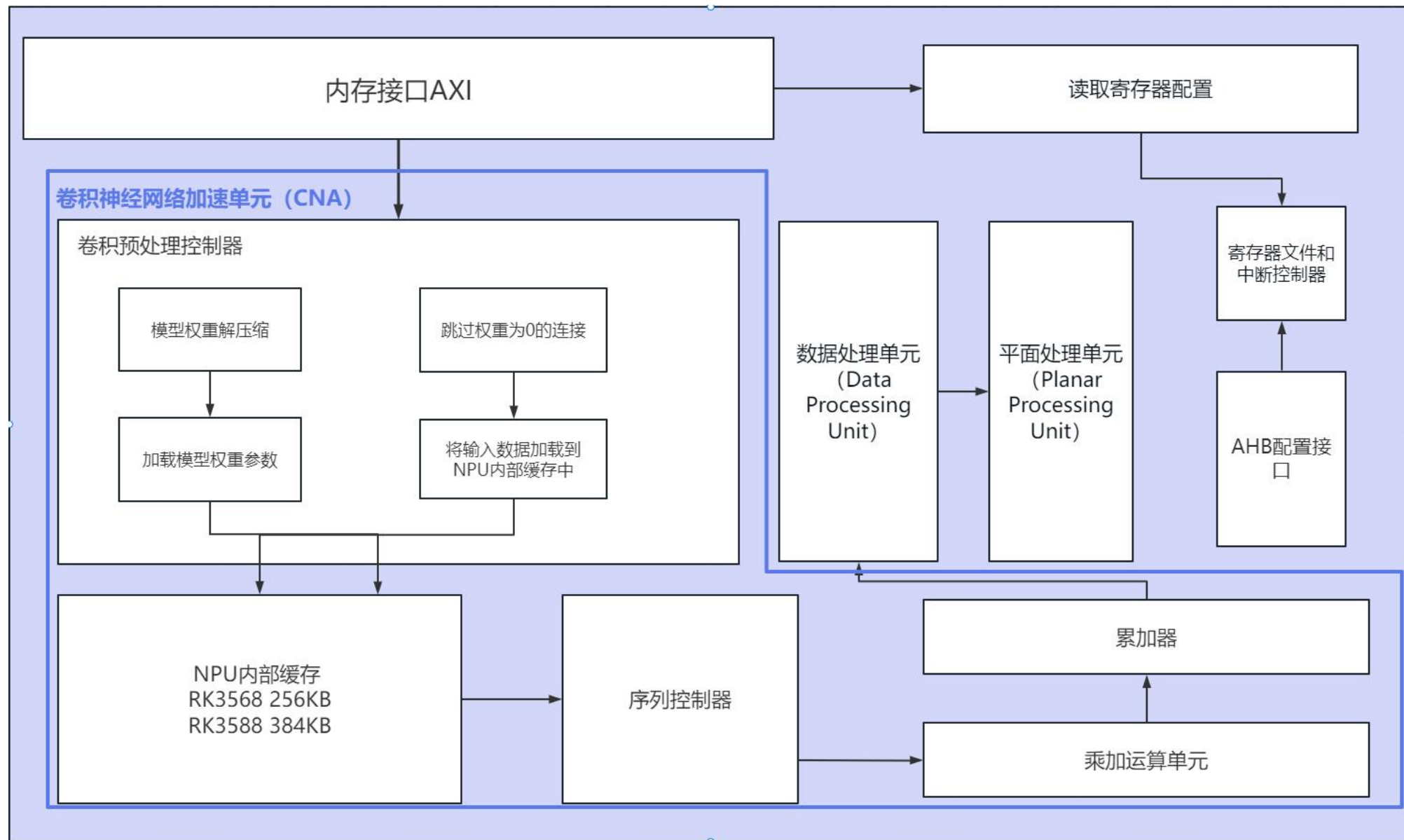
NPU的由来



RKNPU的发展历史



RKNPU单核架构



RKNPU性能计算

NPU算力是指每秒可以处理的运算次数，通常以 $TOPS$ (*Tera Operations Per Second*) 为单位进行衡量。

每个周期的理论峰值性能为 $Perf=MACs*2(ops/cycle)$ 。其中MACs表示每个周期内可以进行的乘加操作次数，而算力指的是运算的次数，所以要将乘加操作分解为一次乘法和一次加法，也就是最后乘以2的由来。

当NPU的频率为 f 时，每秒的理论峰值性能为 $Perf=MACs*2*f(ops/s)$ 。

以RK3588 int 8 数据类型为例进行性能计算演示：

(1) RK3588每个周期可进行 $1024x3$ 个int8 MAC操作

(2) RK3588 的NPU频率为 $1G\ HZ$

(3) 理论峰值性能计算公式为 $Perf=MACs*2*f(ops/s)$

综上RK3588 int 8 理论峰值性能为 $Perf=1024x3x2x1G = 6\ TOPS$

RKNPU应用场景

计算机视觉

图像分类



目标检测



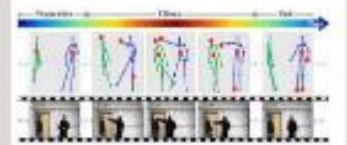
图像分割



人脸识别



行为识别



语音识别

语音识别



语音合成



说话人识别



自然语言处理

情感分析



文本生成



文本分类

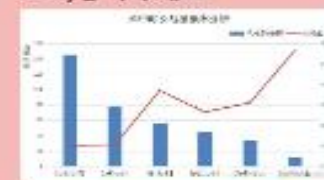


机器翻译



医疗保健

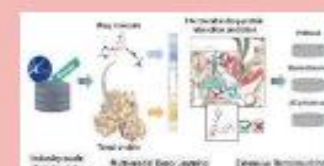
疾病预测



医疗影像分析



药物发现

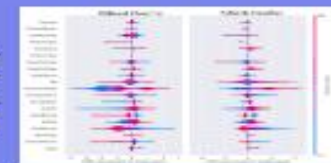


基因组学



金融服务

风险评估



欺诈检测



股票预测



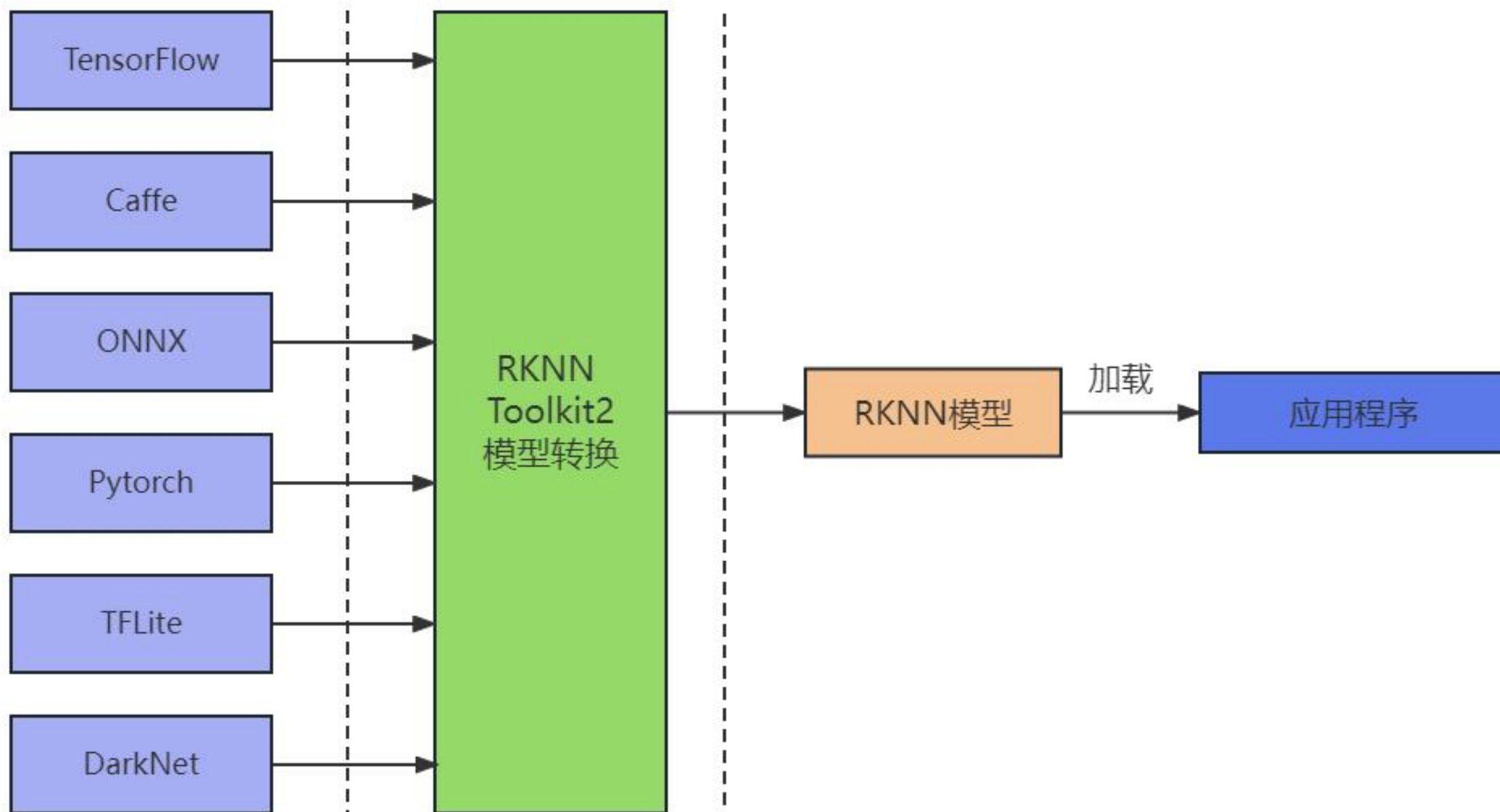
投资组合优化



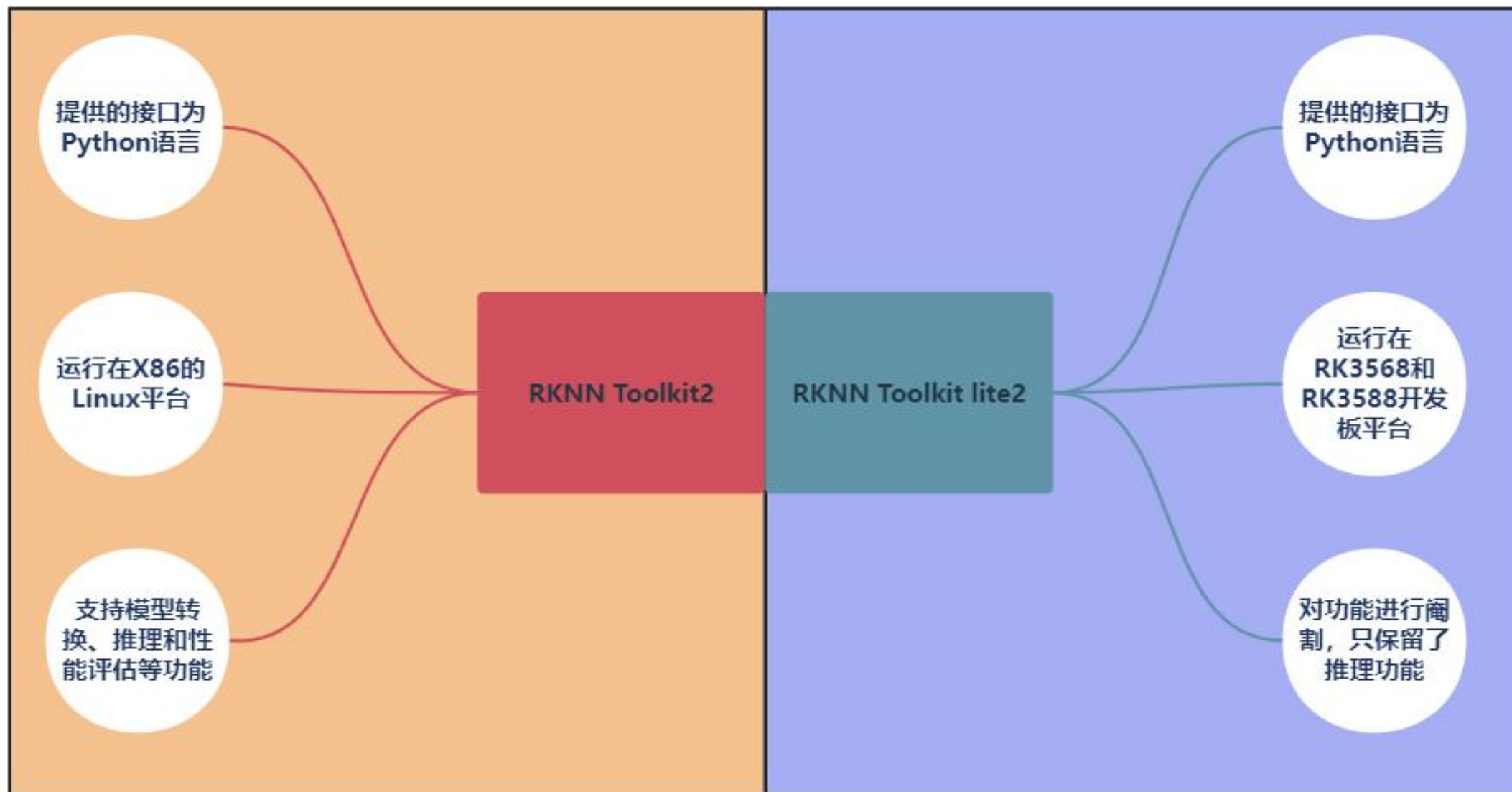
RKNPU 推理软件框架



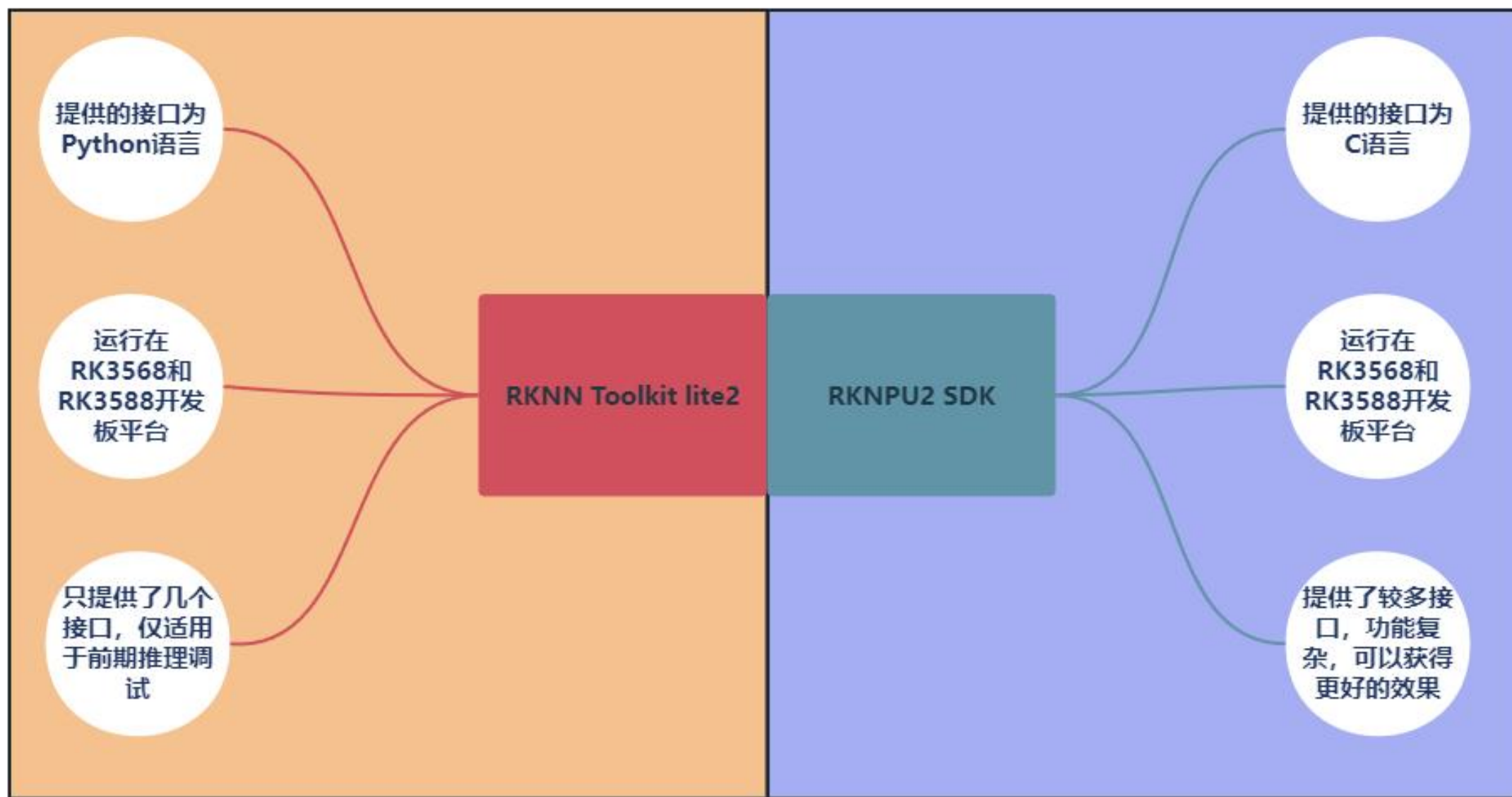
RKNN 模型

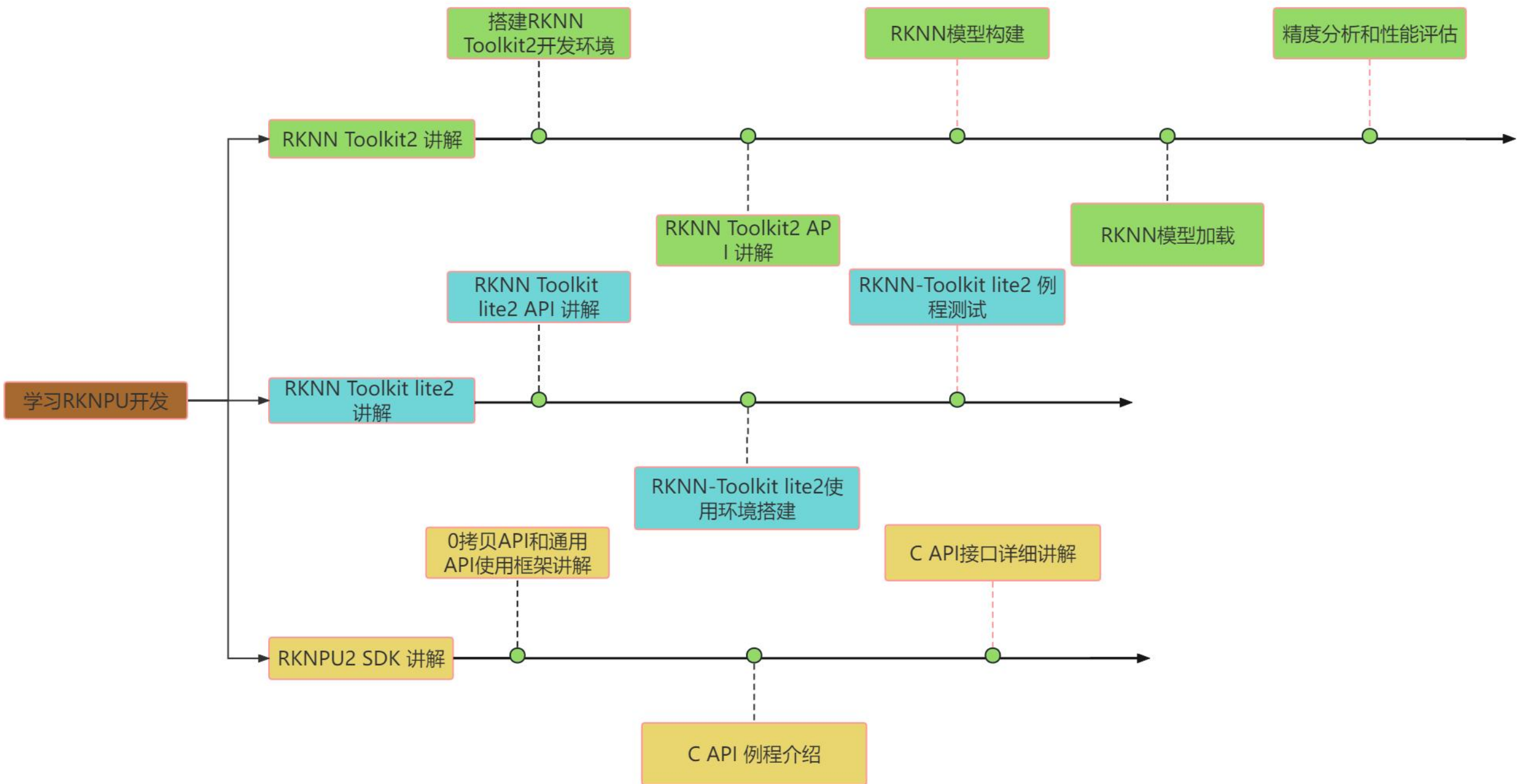


RKNN Toolkit 2和RKNN Toolkit lite2工具对比

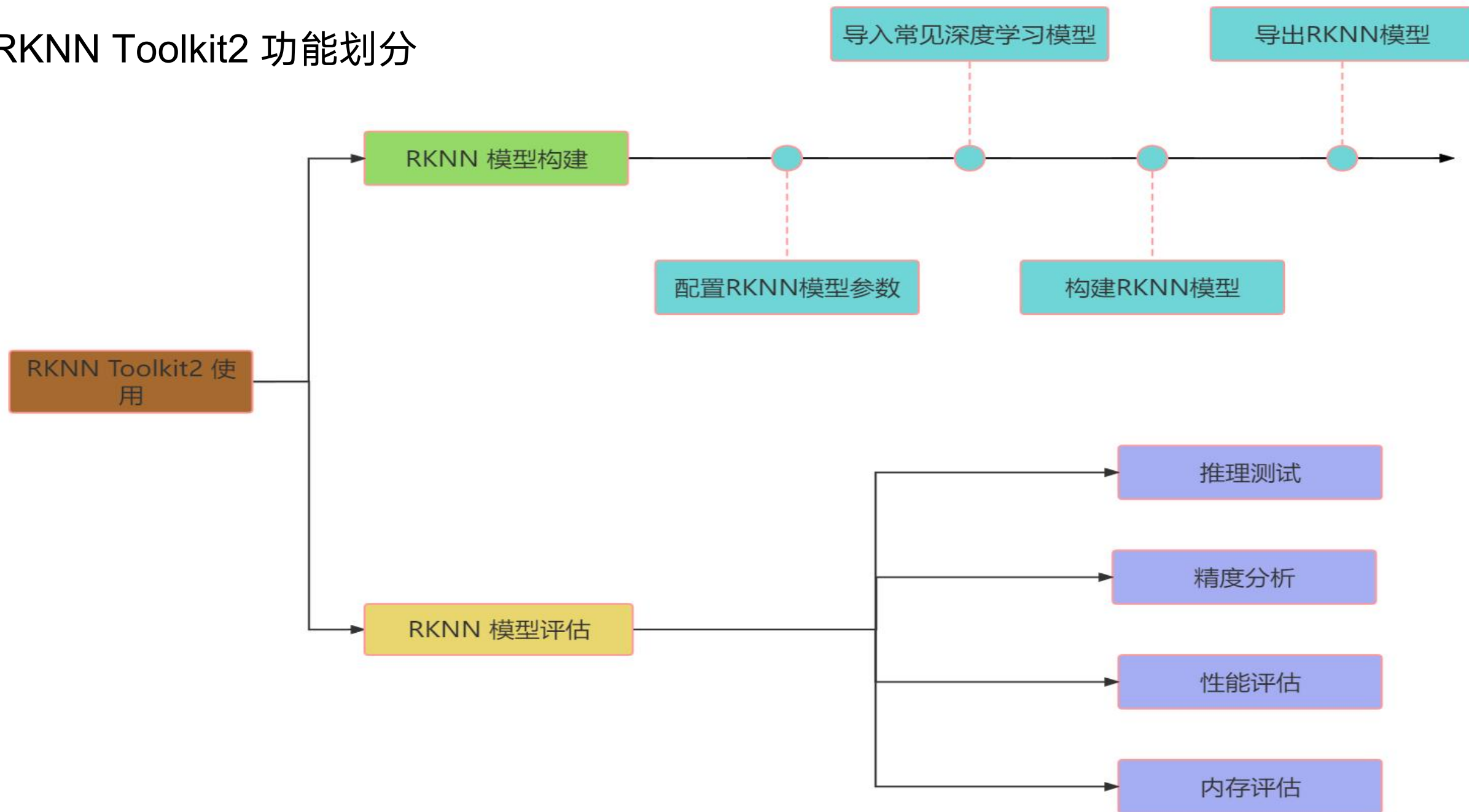


RKNN Toolkit lite2 和RKNPU2 SDK对比

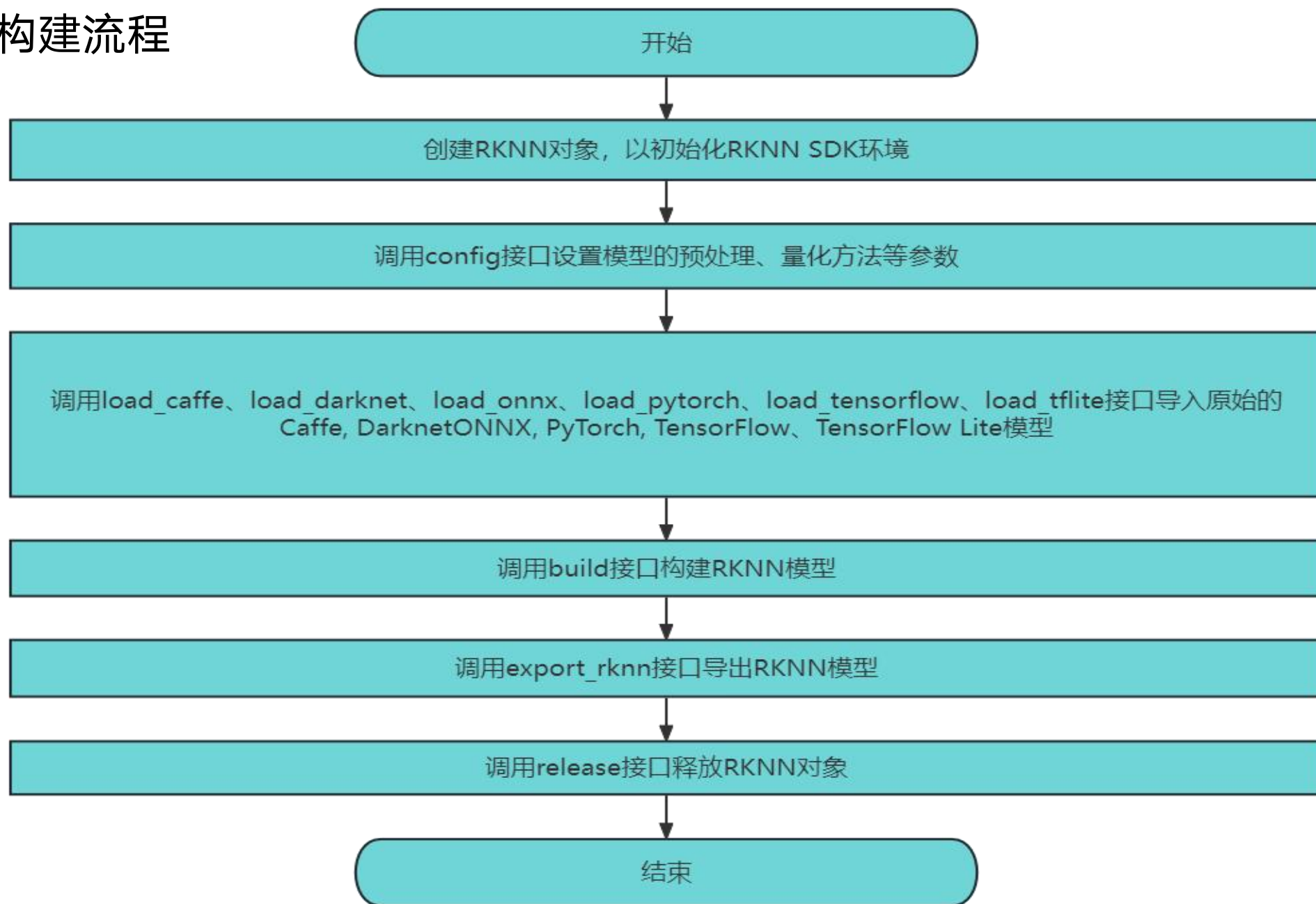




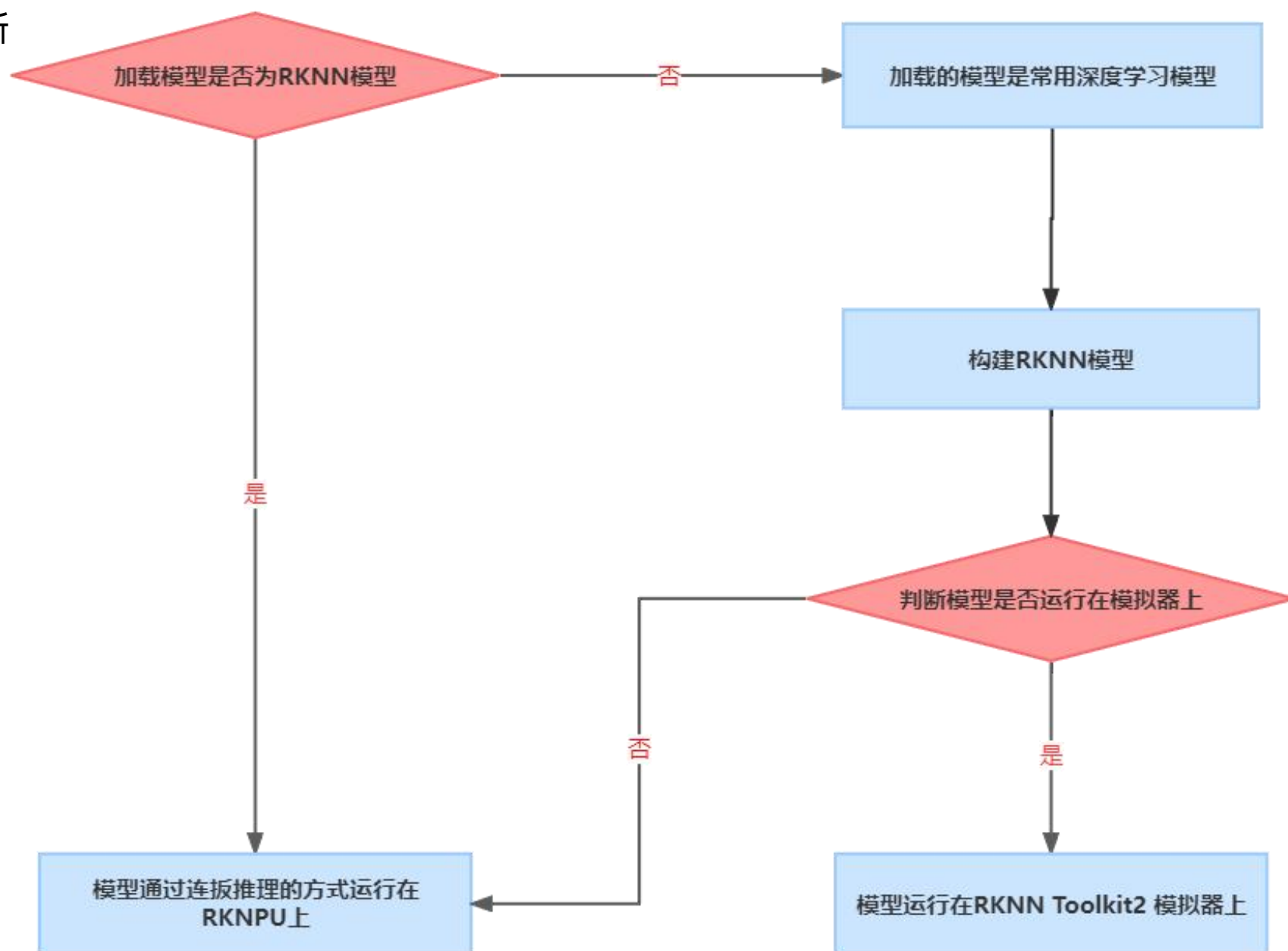
RKNN Toolkit2 功能划分



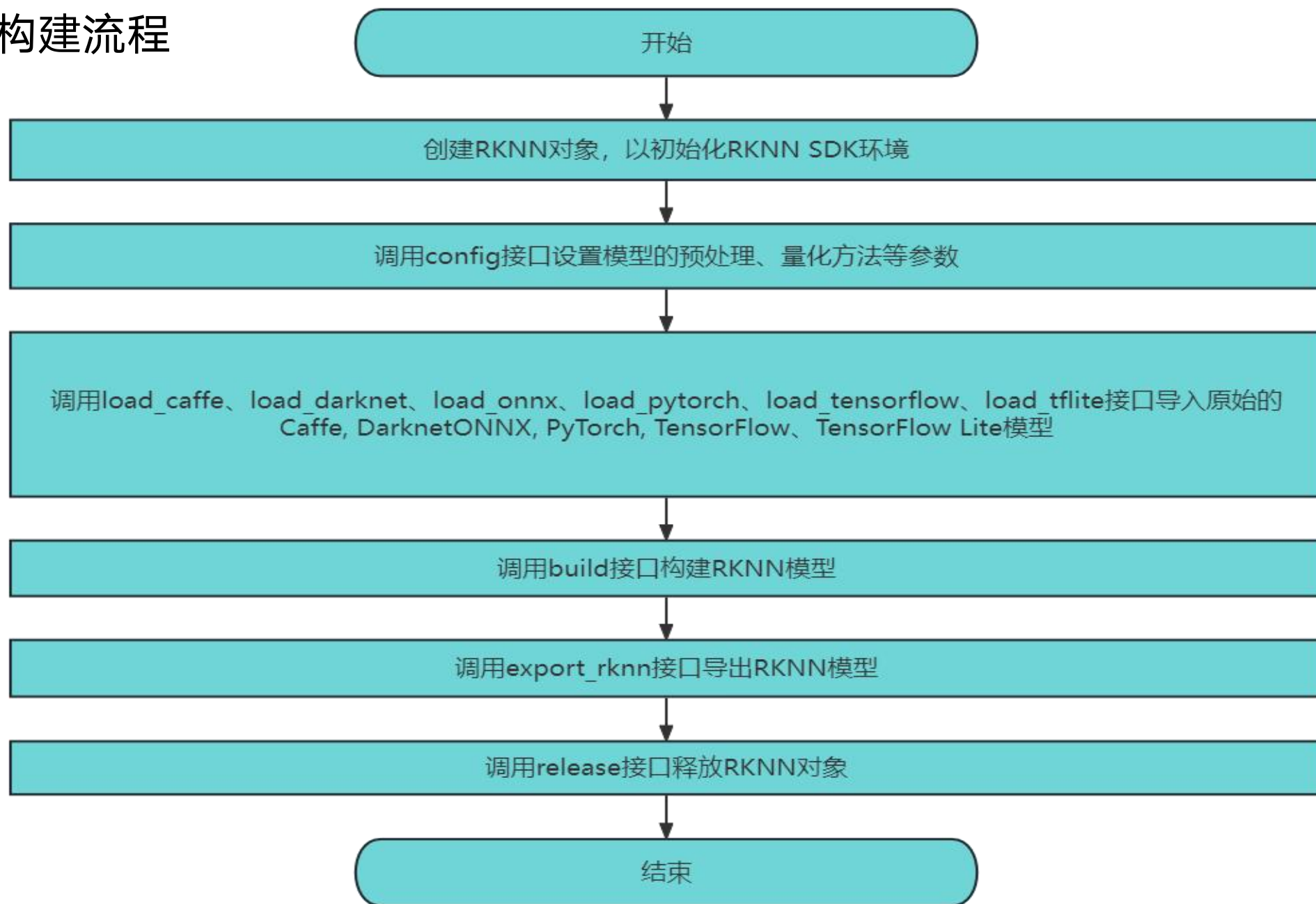
RKNN 模型构建流程



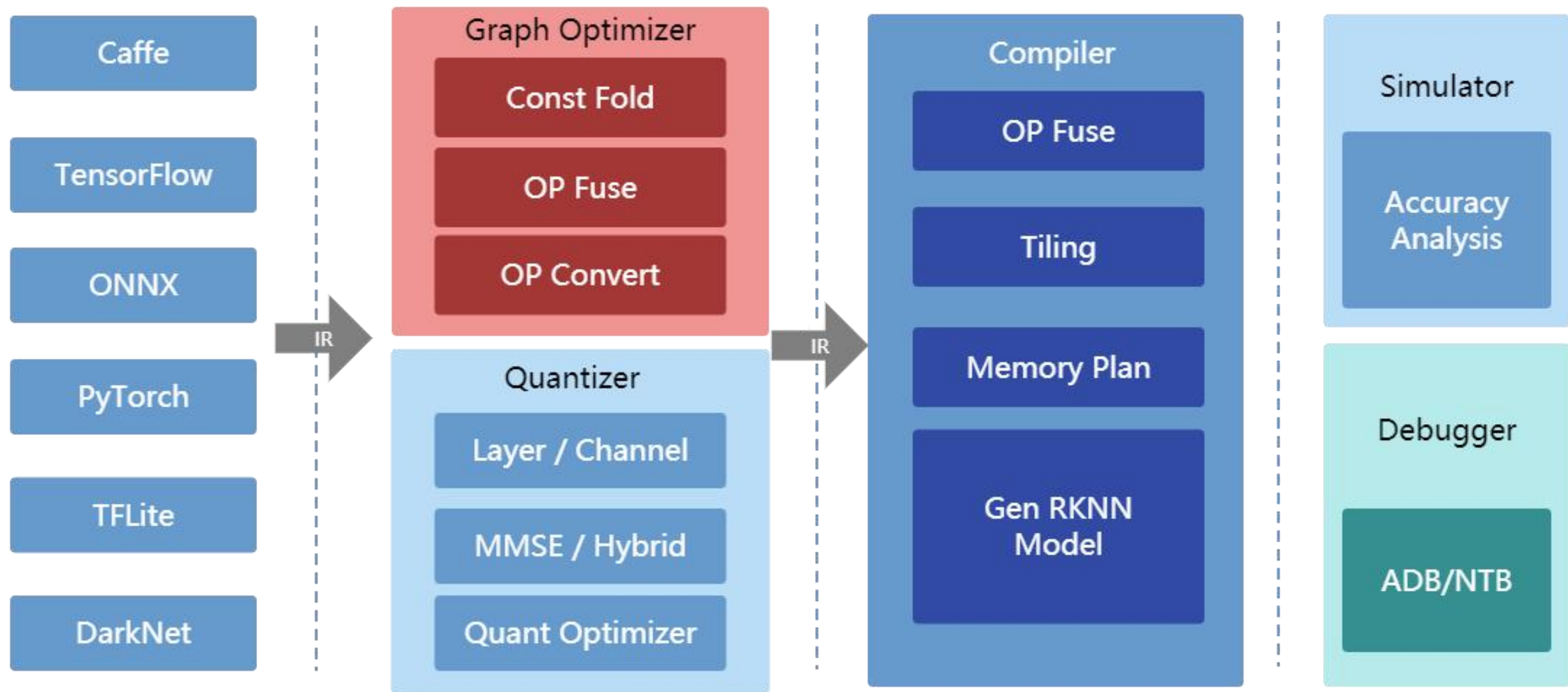
模型运行方式判断



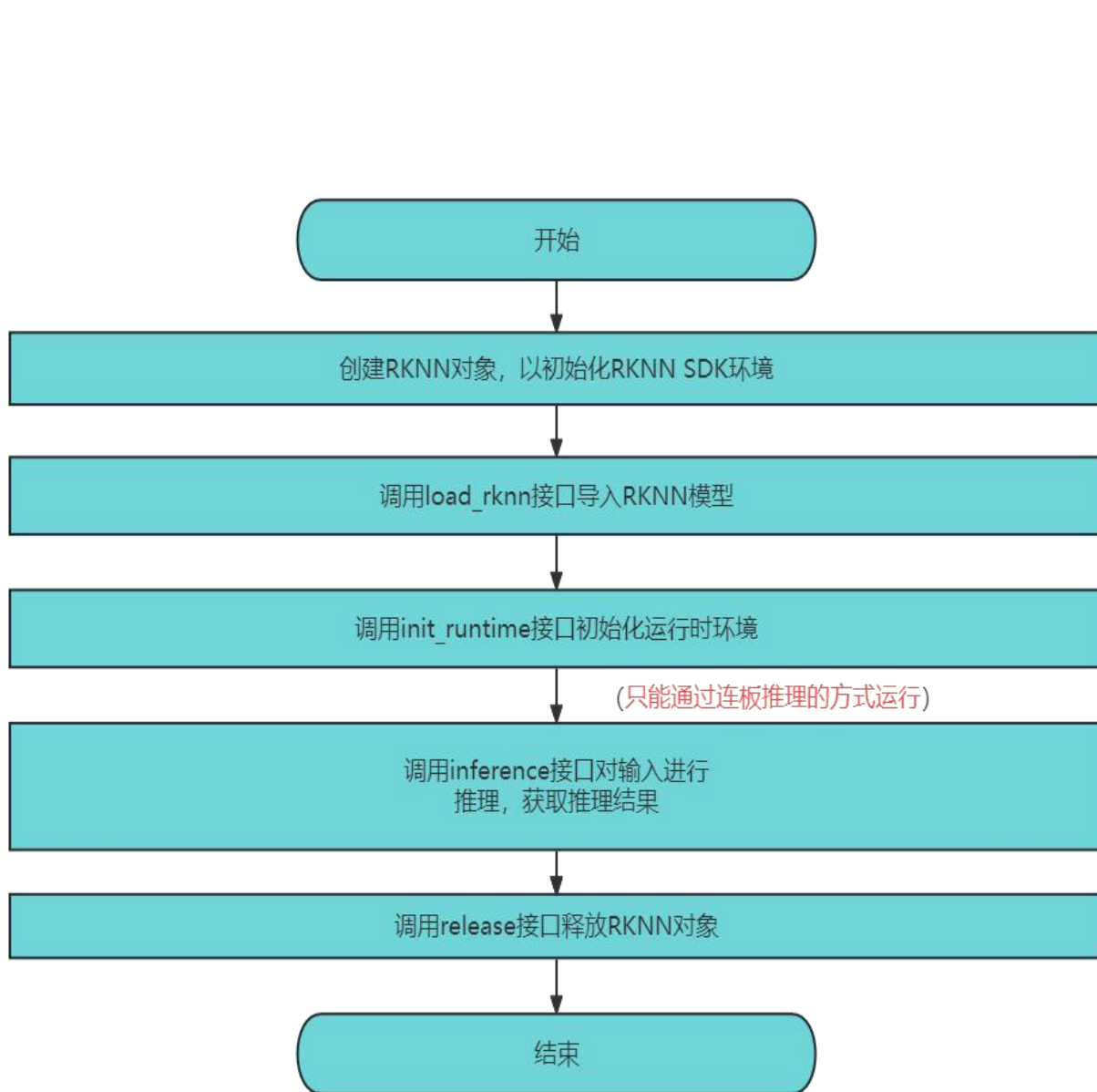
RKNN 模型构建流程



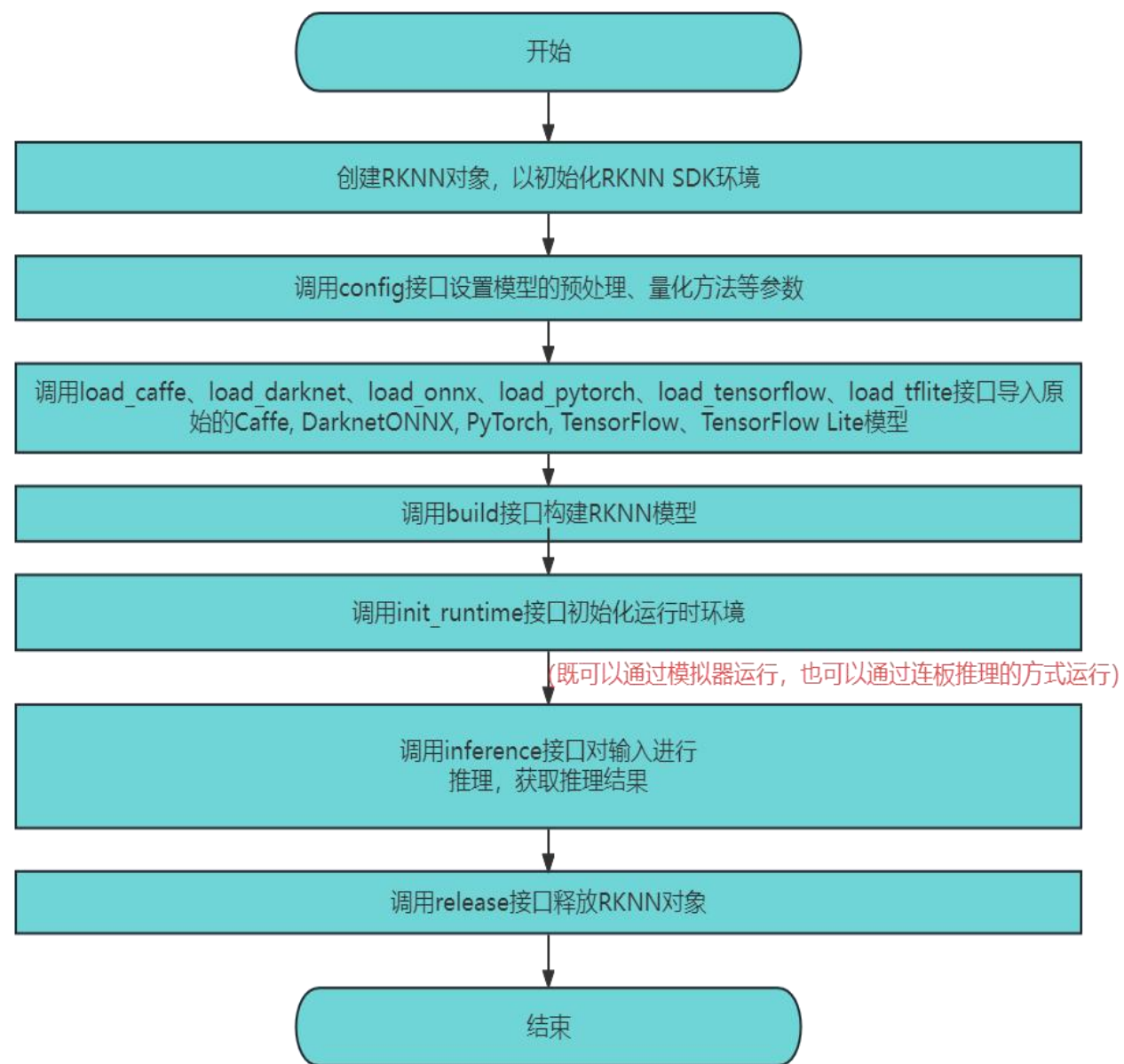
RKNN 模型构建框图



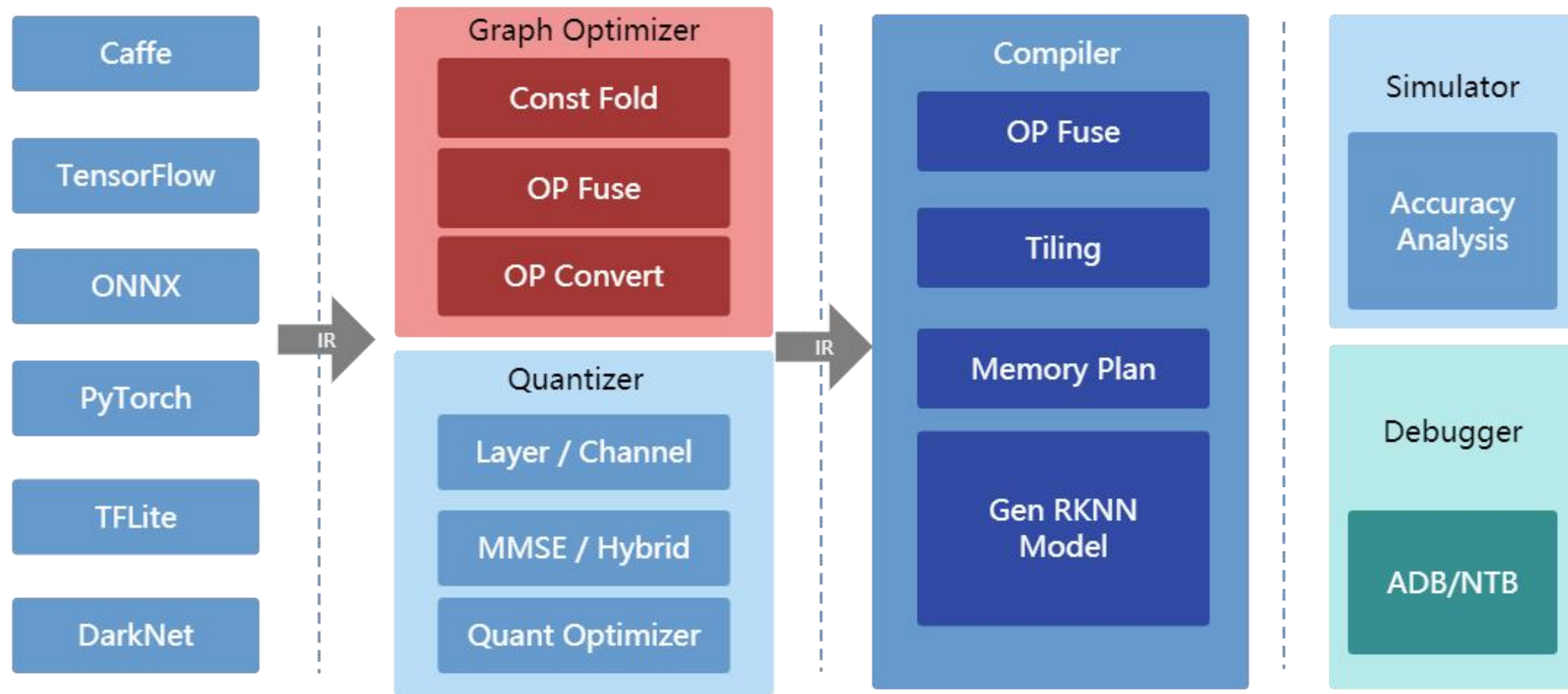
加载RKNN模型进行模型推理



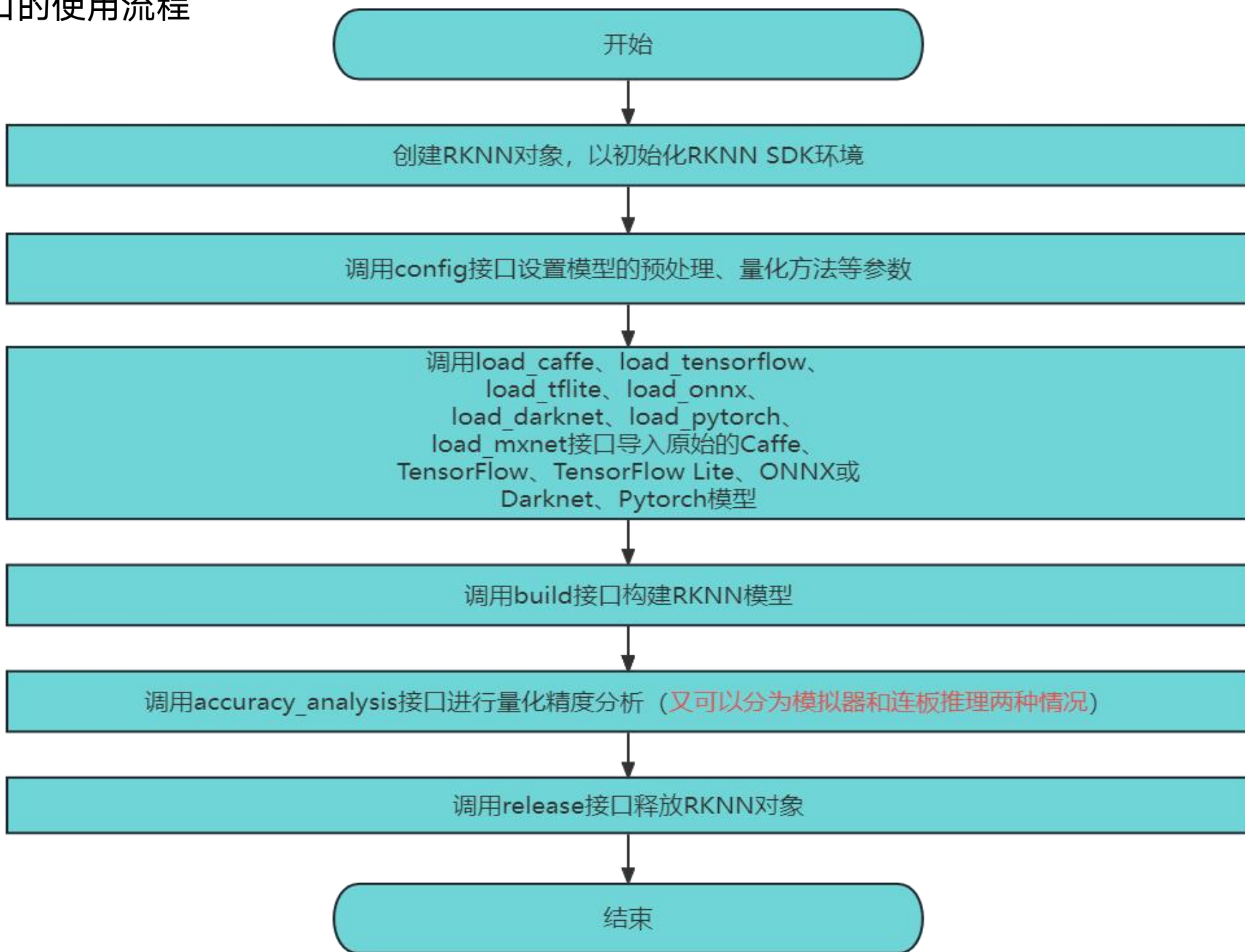
加载非RKNN模型进行模型推理



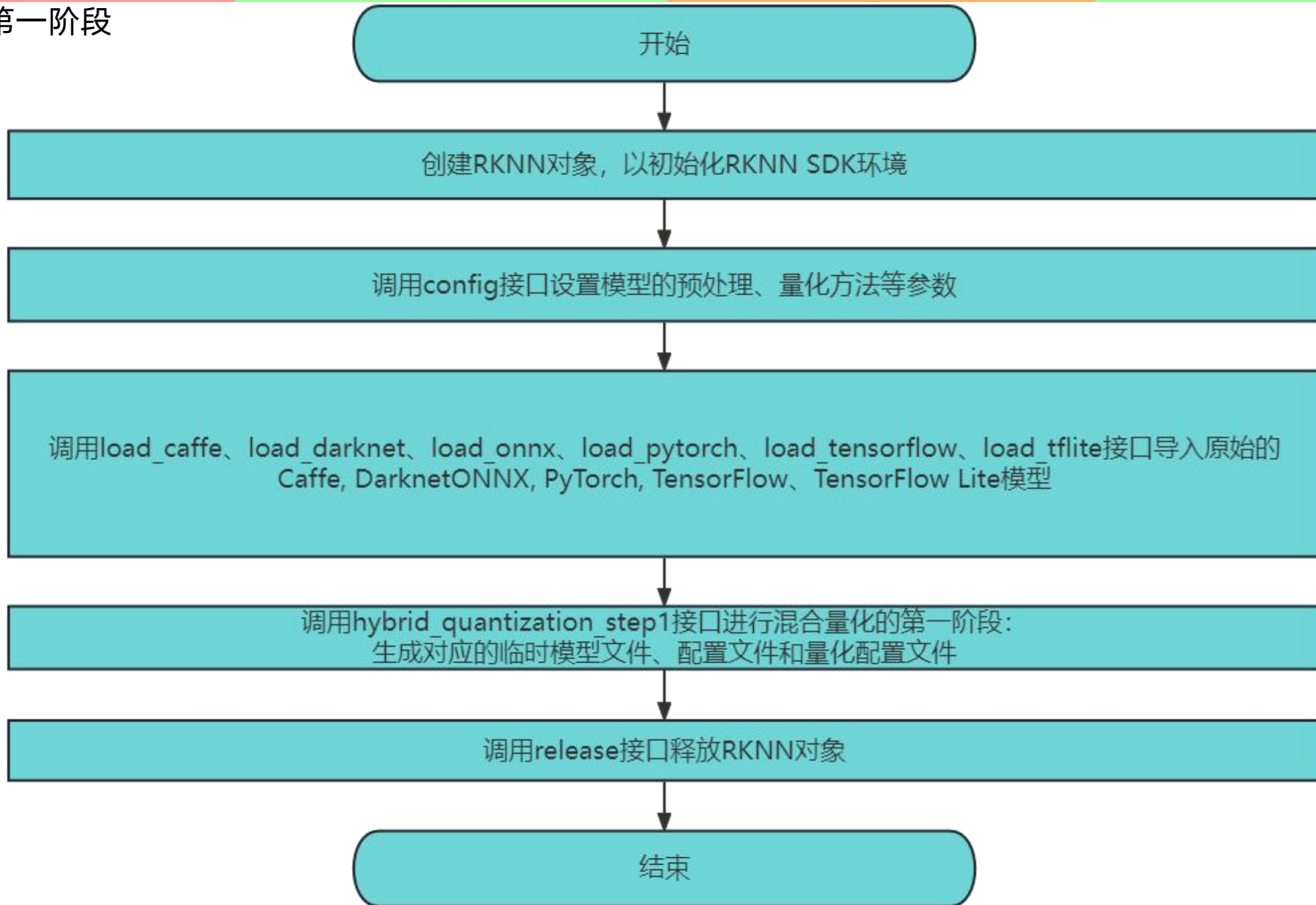
RKNN 模型构建框图



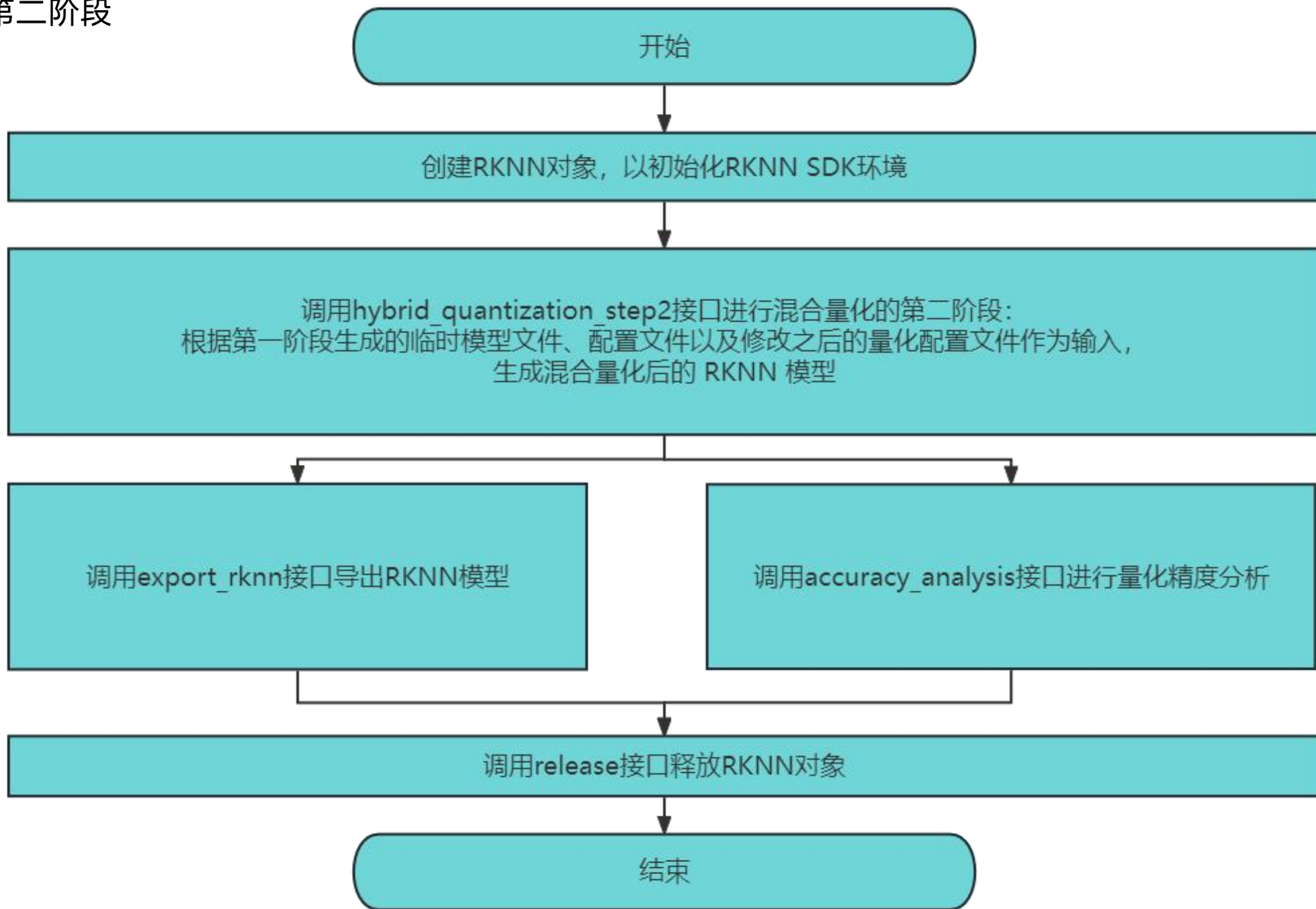
量化精度分析接口的使用流程



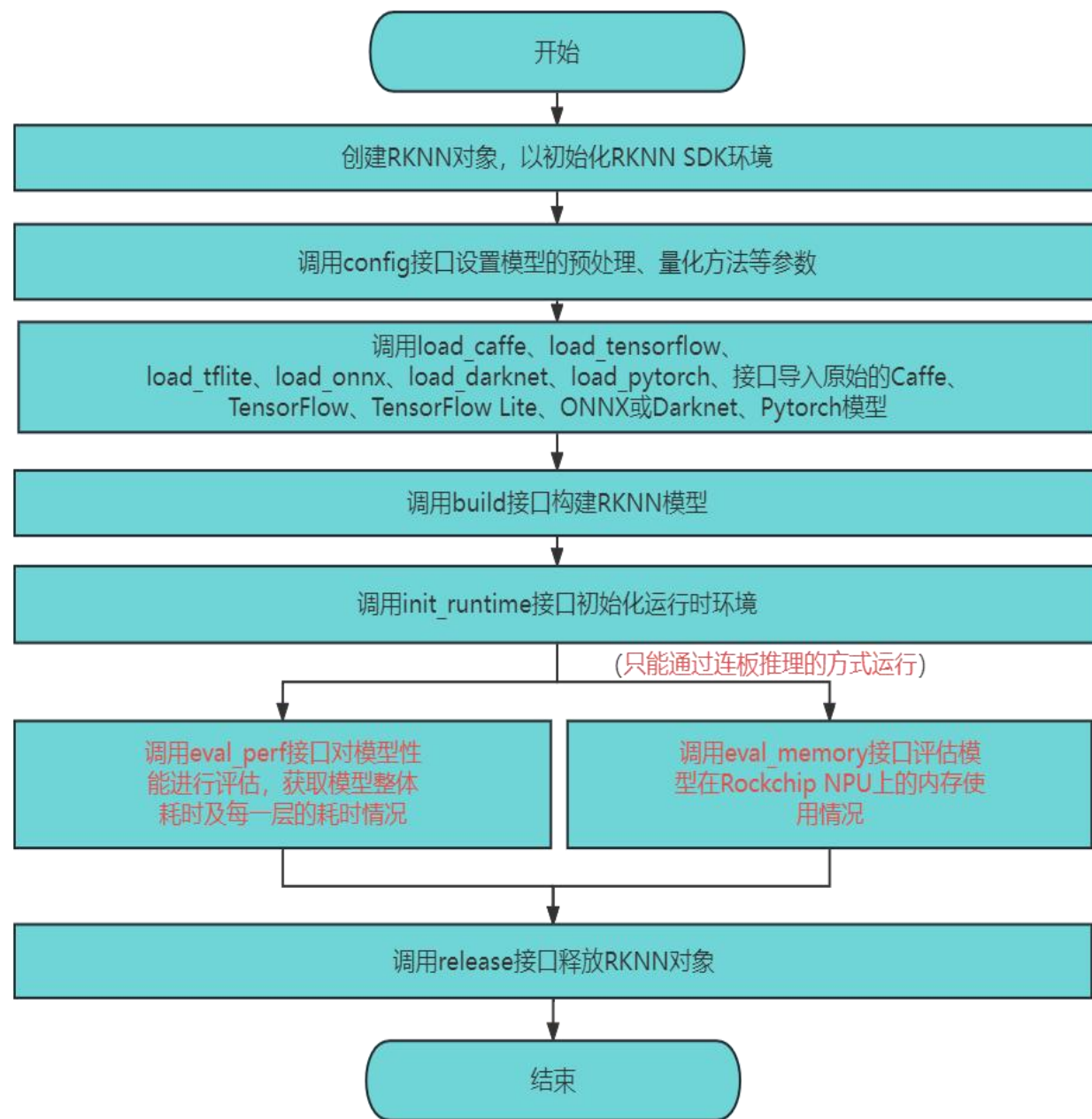
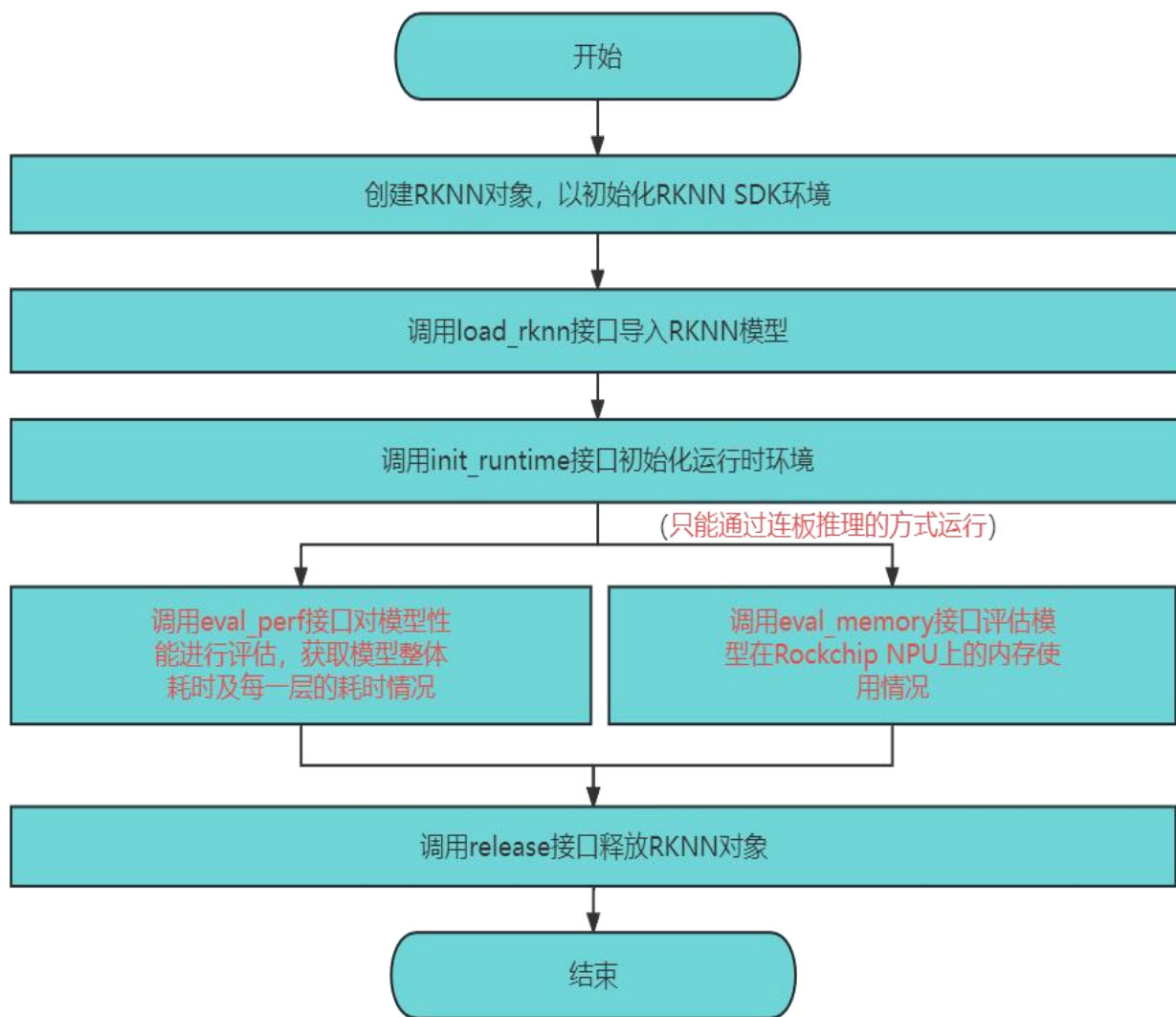
混合量化第一阶段



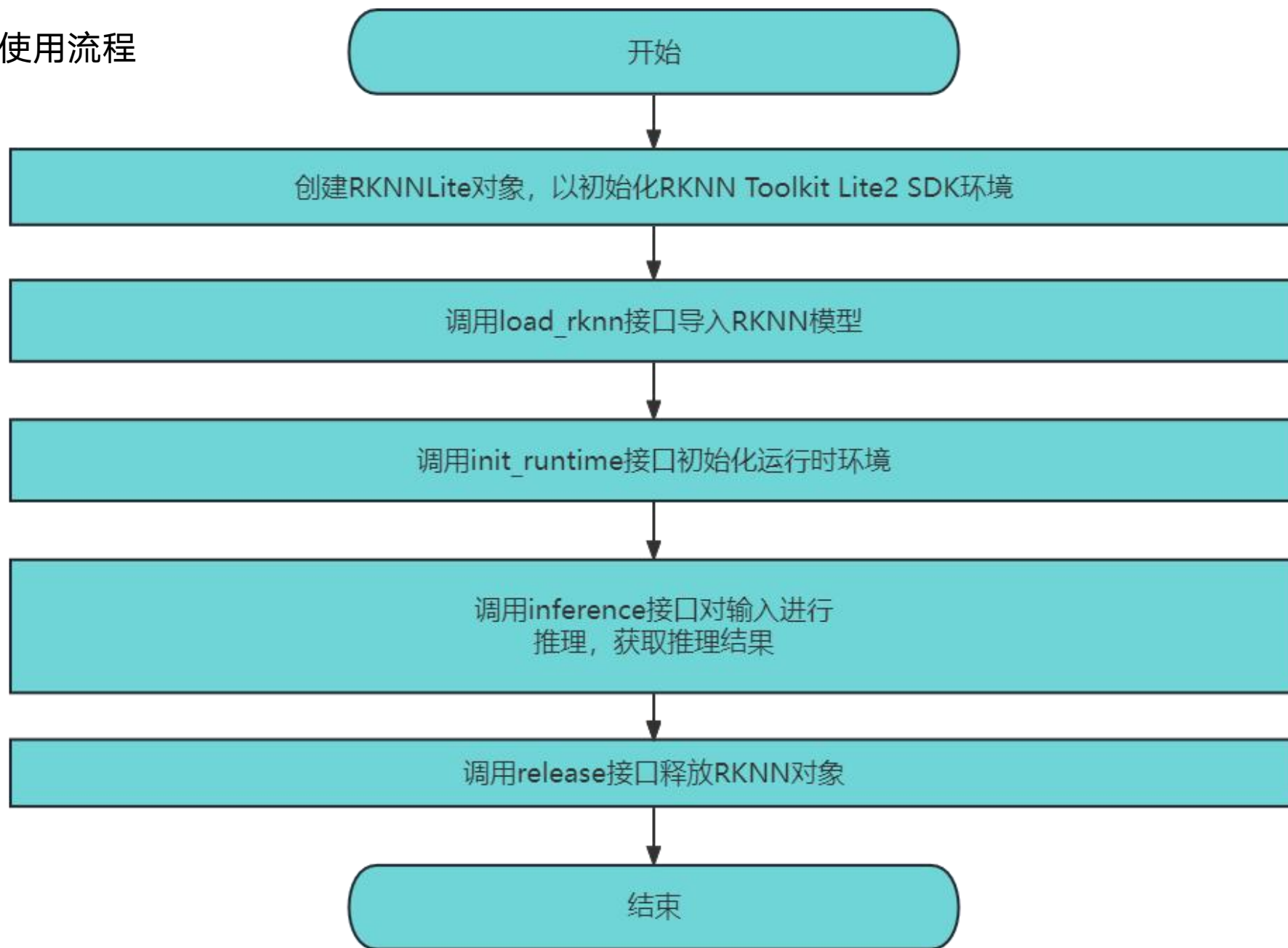
混合量化第二阶段



性能评估和内存评估



RKNN Toolkit lite2 使用流程



RKNPU2 通用API使用流程：



零拷贝API使用流程

