



2020
北京

多媒体开启
MULTIMEDIA BRIDGE
TO A WORLD OF VISION

新视界

视频推理服务 使AI应用更高效

爱奇艺视频推理服务实践

目录

CONTENTS



2020
北京

01 简介

02 现有方案

03 优化方案

04 总结



2020
北京

PART 1

简介

简介



2020
北京

- 周海维
- 爱奇艺深度学习云架构师 (2017-2020)
- 爱奇艺国际部广告平台架构师

- 分享爱奇艺视频推理的工程经验和优化方法
- 重点介绍视频变化服务的工程优化方法
- 不涉及具体算法，任务治理

视频质量增强



2020
北京



片源

增强

SDR2HDR



2020
北京



视频推理服务类型



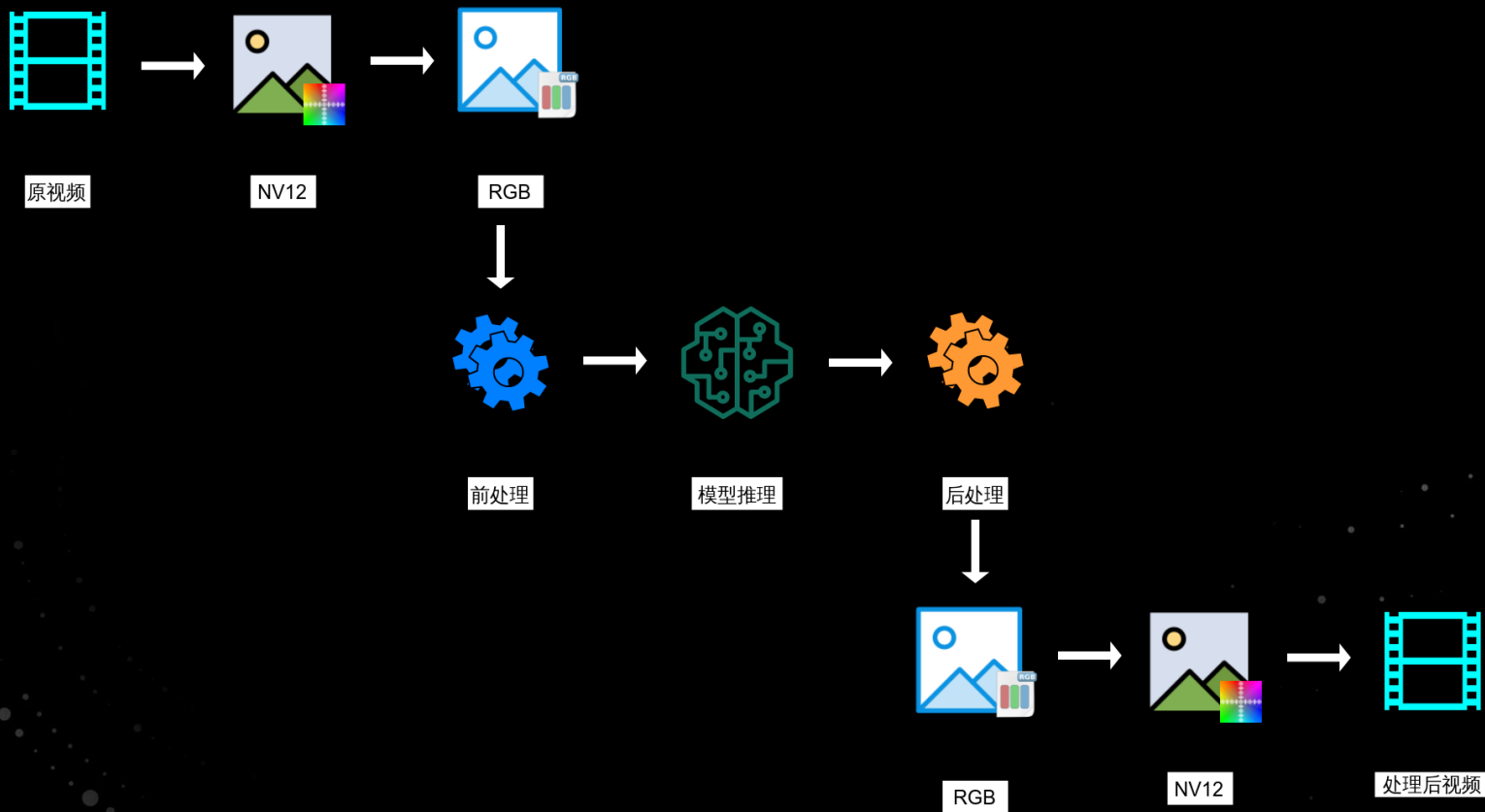
2020
北京

	视频特征提取	视频变化
输入	视频	视频
输出	数值（概率，坐标）	视频
尺寸	缩小	保持或变大
帧率	只留关键帧/采样	保持或增大
计算量	大 可运行在CPU/GPU	巨大 一般运行在GPU
算法	CNN, RCNN	CNN, 传统CV 主观性较强
用途	分类, 物体识别	图像增强, 高帧率, 画面修复

视频变化服务流程



2020
北京





2020
北京

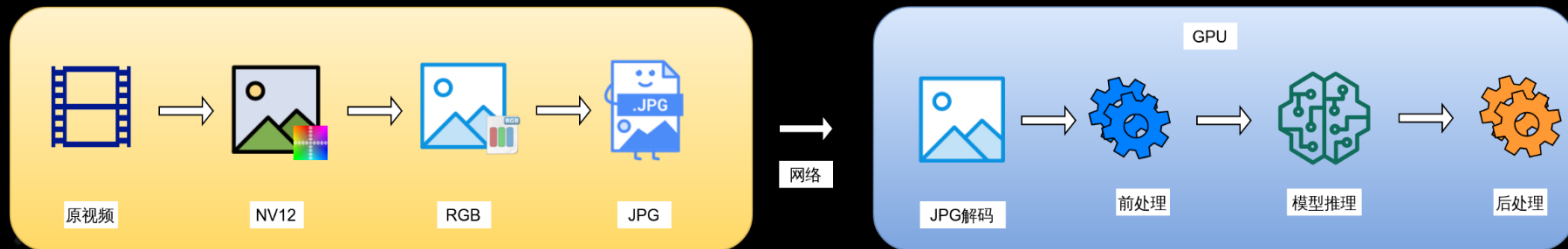
PART 2

现有方案

从 图片推理 演进到 视频推理



2020
北京

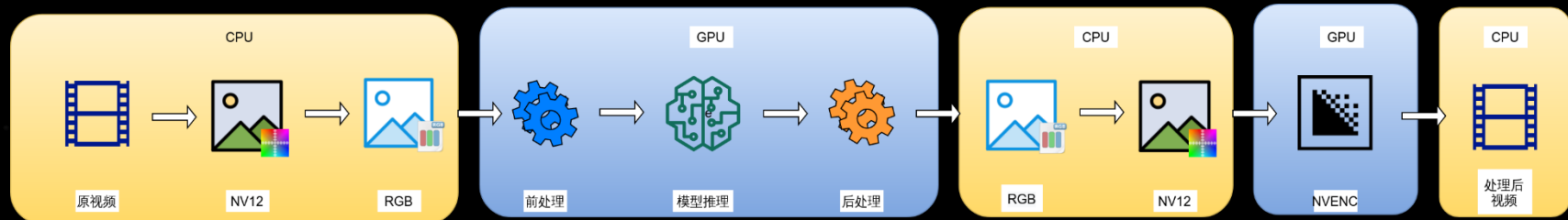


- 单机CPU/GPU比例限制

ffmpeg



2020
北京

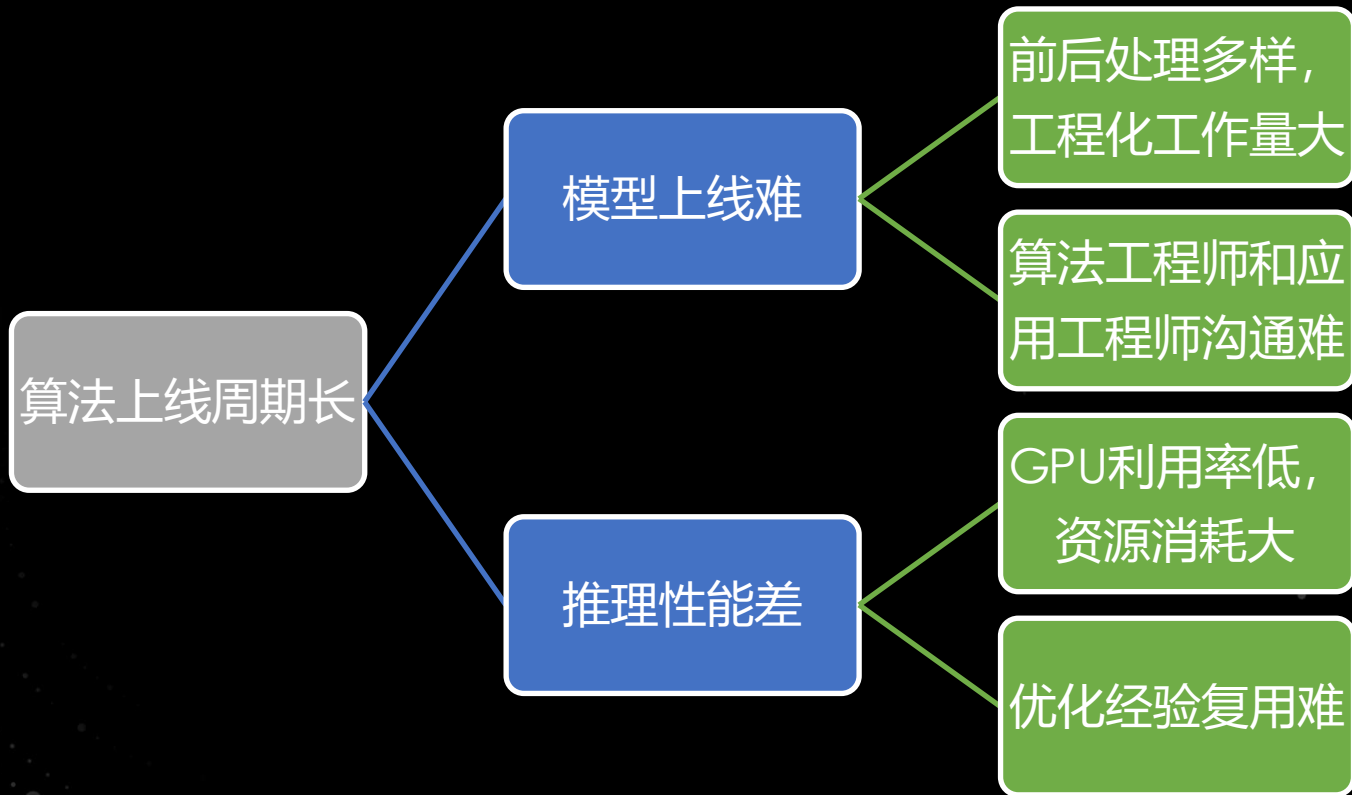


- 1080p RGB 5.93MB
- Pageable Mem <-> GPU 3GB/s
- Pinned Mem <-> GPU 6GB/s
- 不支持batching

已有方案的缺点



2020
北京



GPU常用优化手段



2020
北京

Batching 2x-3x

Nvcodec 2x-10x

Quantization 2-3x

OP fusion 1.2x-10x

Cuda kernel 2-20x



2020
北京

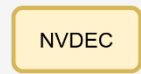
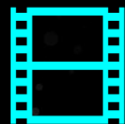
PART 3

优化方案

All in GPU



2020
北京

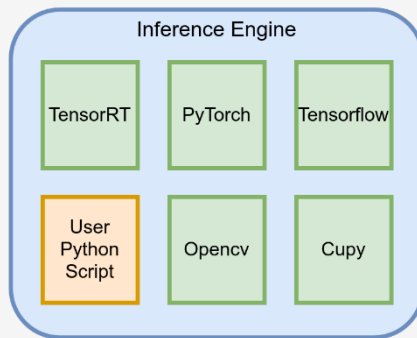


H264
HEVC
HEVC Main10

最大支持
4096*4096



多帧图片



支持
多模型
前后处理脚本



优化后多帧图片



H264
HEVC
HEVC Main10

最大支持
4096*4096

HEVC 1080p
0.86x-4x





2020
北京

Python脚本支持

```
1 def setup(config):
2     # setup model here
3     pass
4
5 def warmup(_):
6     # warmup model
7     pass
8
9 def configure(stream_config):
10    # configure per stream
11    pass
12
13 def process(data, profiling=False):
14    # batching process frames
15    pass
```

```
1 import numpy as np
2 import cupy as cp
3 from PIL import Image
4
5 setup('{"model":"/data/haiwei/resnet_v2_fp32_savedmodel_NCHW/"})
6
7 warmup()
8
9 configure('{"user-param" : "test"}')
10
11 im = Image.open('test.jpg')
12
13 data = {
14     'height' : im.height,
15     'width' : im.width,
16     'format' : 'RGB',
17     'data': [ cp.asarray(np.array(im)).data.ptr for i in range(16)]
18 }
19
20 process(data)
```


Python脚本支持



2020
北京

算法工程师更容易接收Python



开发效率比C++高，运行效率没有显著变化



兼容性比C++动态库好

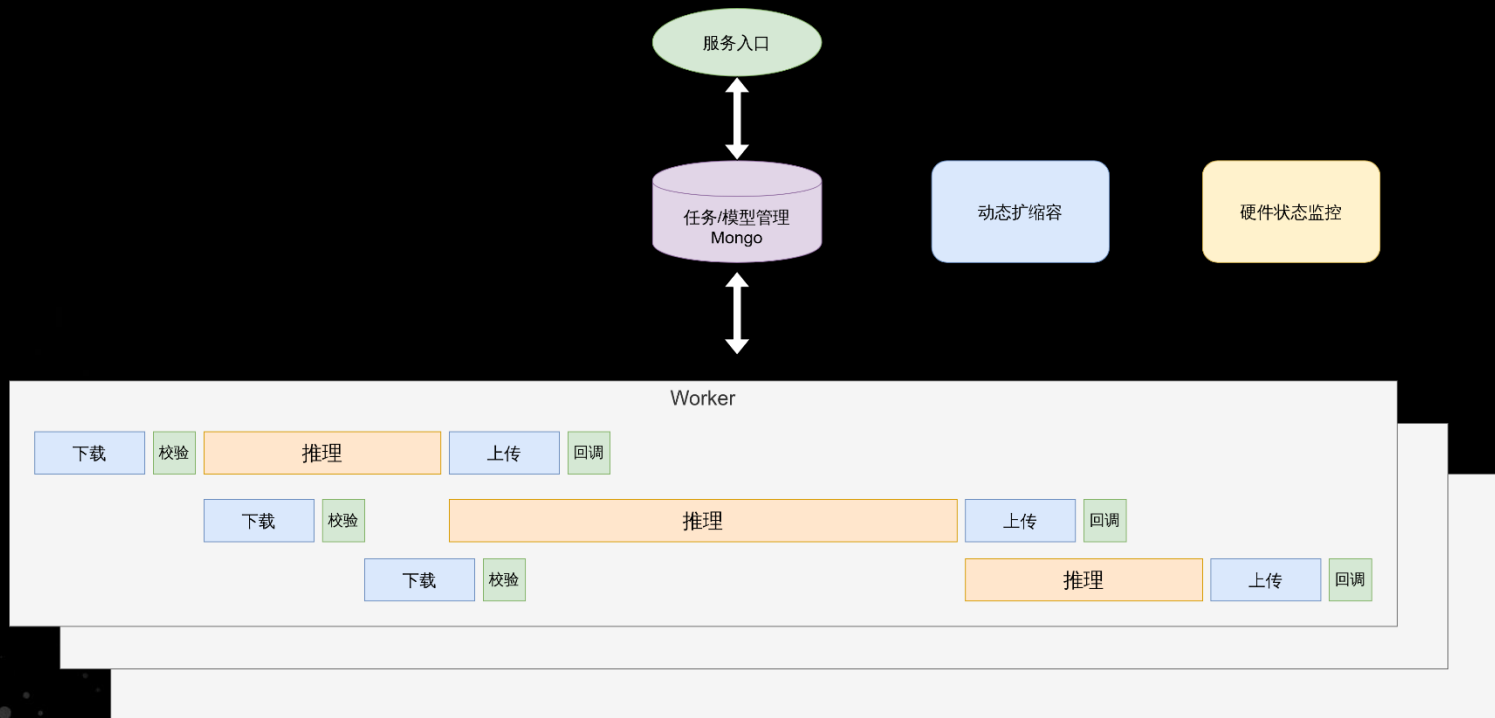


平台工程师需求降低，上线周期更短

任务流水线



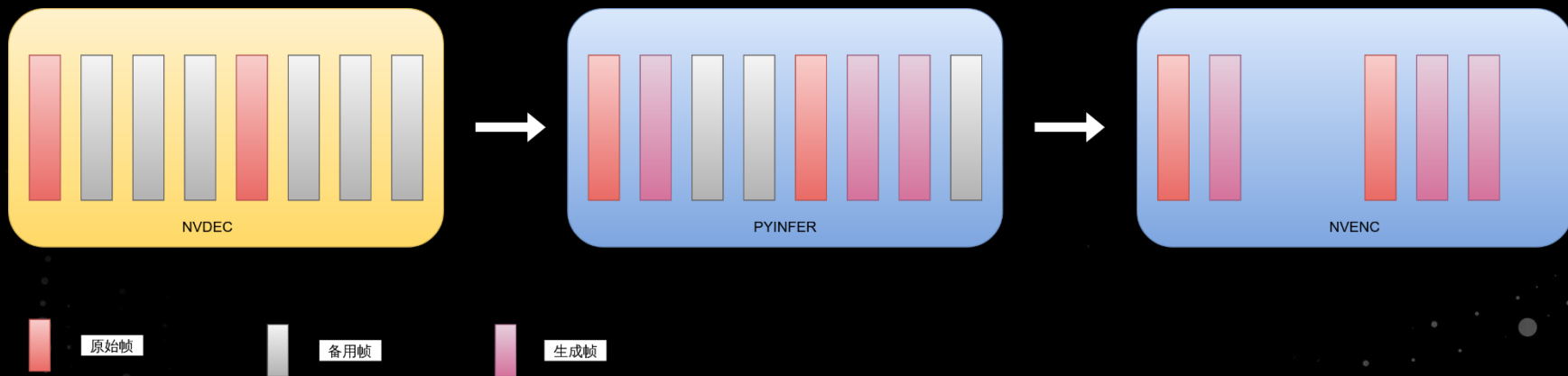
2020
北京



高帧率变换服务



2020
北京

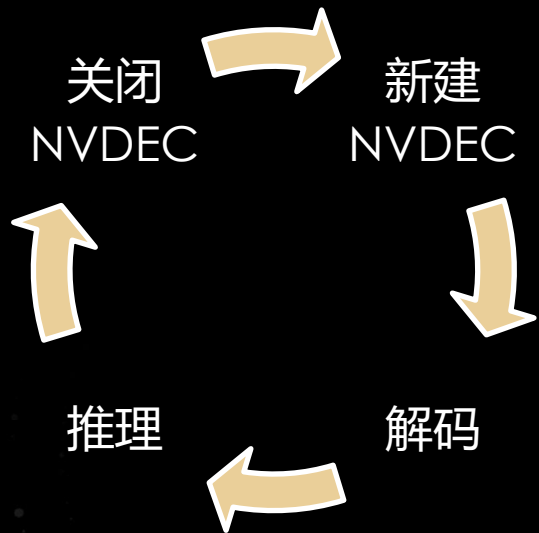


- 解码时按比例添加备用帧
- 编码时保留原始帧和生成帧，跳过备用帧

短视频分类服务



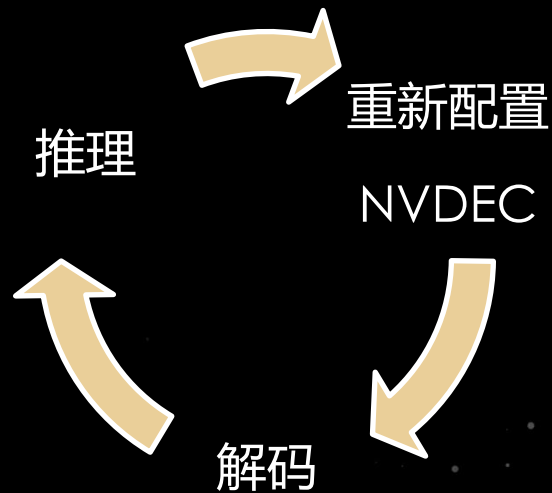
2020
北京



110ms



=>



15ms

模型效果偏差问题



2020
北京

现象

- 少数推理结果差异较大

原因

- 图像伸缩算法不一致
- 色彩空间转换取整误差

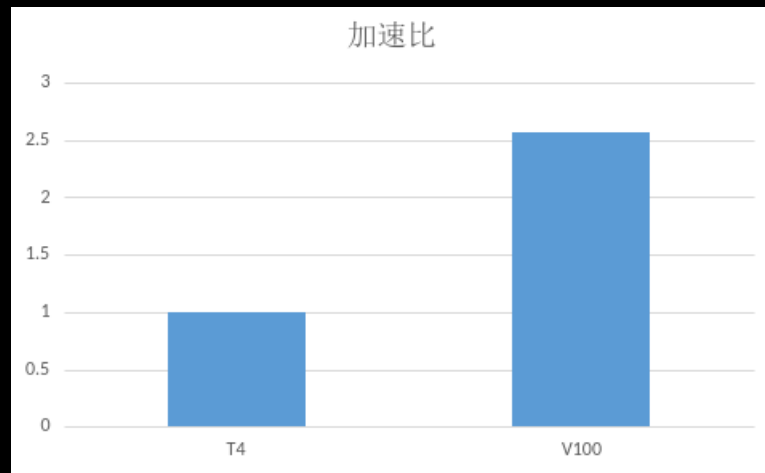
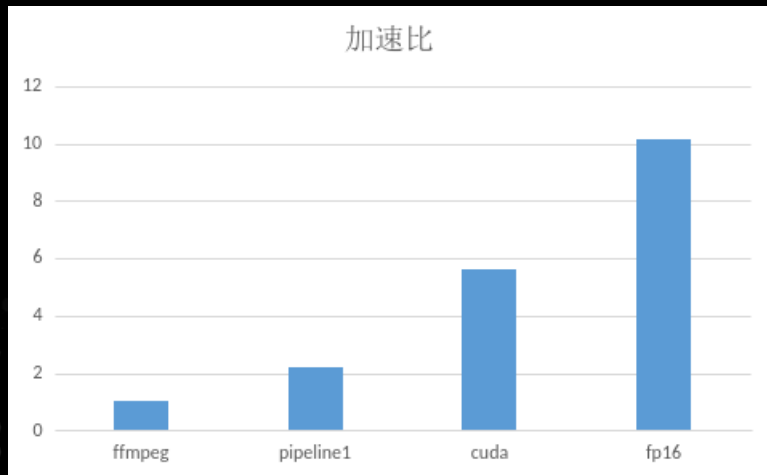
方案

- 使用同一套预处理操作，重新训练模型

优化效果



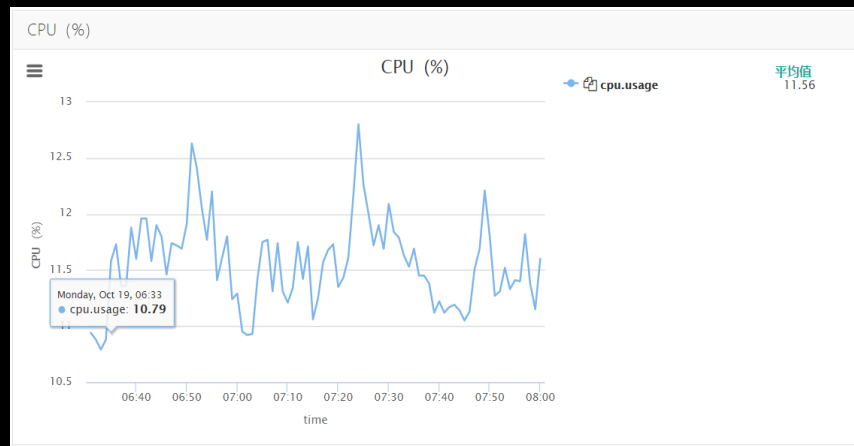
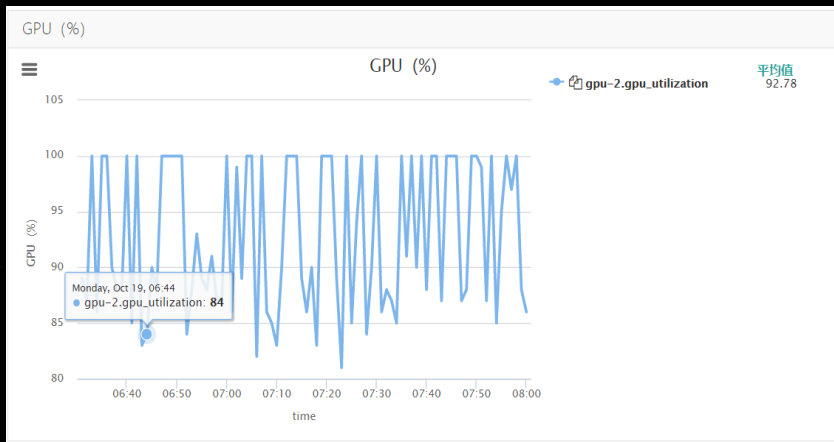
2020
北京



资源利用率



2020
北京





2020
北京

PART 4

总结

应用经验



2020
北京



总结



2020
北京

集成的优化措施

- 视频推理过程集中在GPU上完成，减少CPU/GPU数据搬运开销
- 使用任务流水线方式，并行CPU/GPU操作
- 减少重复优化

集成Python引擎

- 满足灵活算法需求
- 降低使用门槛

统一的开发测试环境

- 算法工程师专注于算法处理逻辑
- 平台工程师专注于性能优化



2020
北京

多媒体开启
MULTIMEDIA BRIDGE
TO A WORLD OF VISION

新视界

Thank you

