

# 基于FPGA加速的 工业缺陷检测系统

A组 CICC1577 冰糖葫芦儿

成员: 纪佳骏 张峰 宋德浩





# 目录

1 团队介绍与研发情况

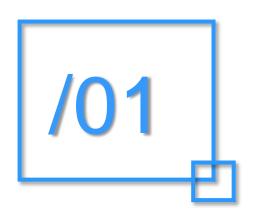
2 效果演示与性能指标

03 关键技术分析

04 总结与补充







## 团队介绍与研发情况

- -团队介绍
- ·项目进展



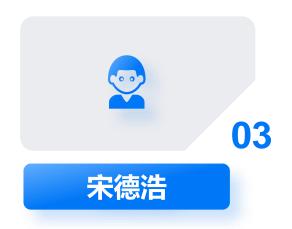
## 团队介绍











专业: 电子信息工程

负责: Linux工程优化

专业:通信工程

负责: FPGA工程优化

专业: 电子信息工程

负责:模型优化



## 项目进展

240ms

3月

Demo整合







112ms

4月

SDK优化、FPGA工程 优化、频率提升、 模型量化训练 5月

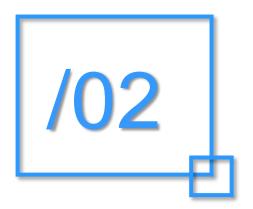
重写DVP、VGA IP, 时序约束、逻辑固化、 剪枝量化联合训练、 数据重排PL加速可行性 验证 6、7月

剪枝算法优化、引入PL侧绘制预测框IP、PL侧保电size加速IP、DVP\_VGA连续传输、增量编译、工程时序优化等









## 效果演示与性能指标



## 效果演示—单次推理 (DEMO测试)

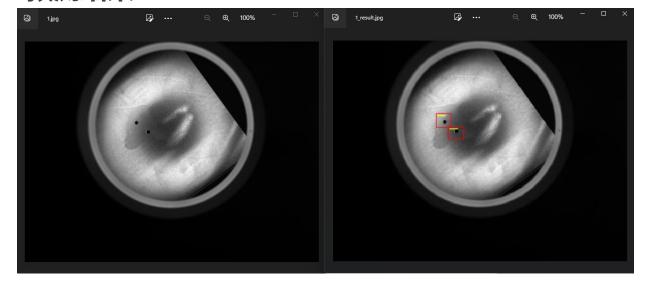




#### 推理速率

```
input organize time:10.588405ms
fpga time:37.580349ms
output organize time:3.950532ms
 ----graph end-----
 ----graph begin-----
input_organize_time:9.738510ms
fpga_time:37.172298ms
output_organize_time:3.657476ms
-----graph end-----
success detection, image size: 640, 480, detect object: zhen_kong, score: 0.996985, location: x
=236, y=157, width=23, height=25
success detection, image size: 640, 480, detect object: zhen_kong, score: 0.784001, location: >
=267, y=186, width=35, height=30
runtime device: armv8
precision: int8/float3
num threads: 2
------ Model info
Model name: new1ch 3.nb
------Perf info -----
Total number of predicted data: 1 and total time spent(s): 77
preproce time(ms): 12.7588, inference time(ms): 64.1086, postprocess time(ms): 1.05696
fpga release
root@awcloud:/opt/paddle frame#
```

#### 预测结果





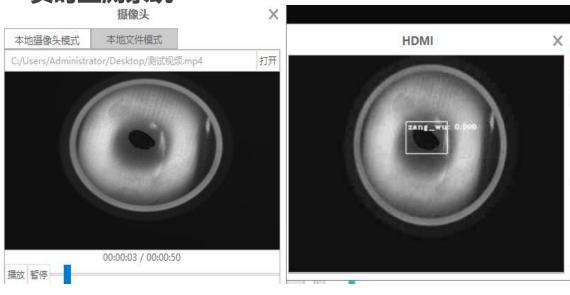
该程序使用SD卡中图片进行推理,可以看到推理时间为64.1ms。并且我们对预处理时间进行了优化,达到了12.75ms。推理结果写入SD卡,如右图所示,可以看到成功检测出了缺陷(这里我们使用了较难监测出的"针孔"缺陷,以体现我们模型的泛化性)。

## 效果演示





#### 实时监测系统



- 重构DVP\_VGA IP实现连续传输,使视频的推流无需 PS控制即可自动完成。
- 重构PL\_Plot模块,完成PL侧绘制预测框,PS侧通过 Avalon\_lw桥控制状态寄存器,实现预测框的实时显示。

#### 进程演示

使用官方提供的LOG参考范式,完成推理结果的输出与时间统计。可以看到推理结果的刷新时间为69ms。



### 性能指标





#### a.Demo单次推理时间:

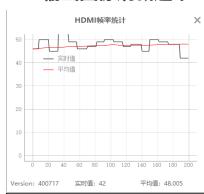
```
-----graph begin-----
input organize time: 10.588405ms
fpga time:37.580349ms
output_organize_time:3.950532ms
 -----graph end-----
 ----graph begin-----
input organize time: 9.738510ms
fpga time:37.172298ms
output_organize_time:3.657476ms
----graph end-----
predictor run:64.015ms
success detection, image size: 640, 480, detect object: zhen_kong, score: 0.996985, location: x
=236, y=157, width=23, height=25
success detection, image size: 640, 480, detect object: zhen kong, score: 0.784001, location: x
=267, y=186, width=35, height=30
----- Config info -----
runtime device: armv8
precision: int8/float3
num threads: 2
----- Data info ------
batch size: 1
Model name: new1ch 3.nb
-----Perf info -----
Total number of predicted data: 1 and total time spent(s): 77
preproce_time(ms): 12.7588, inference_time(ms): 64.1086, postprocess_time(ms): 1.05696 fpga release
root@awcloud:/opt/paddle frame#
```

Inference\_time: 64. 108ms



创新点:我们的单次推理时间与推理刷新时间相差仅5ms!!

#### b1.输出图像刷新速率: 平均值: 48帧



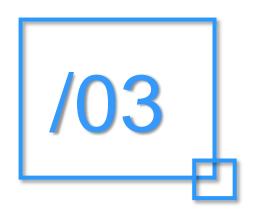
HDMI refresh time: 20.808ms

#### b2.推理结果刷新速率:

Inference\_refresh\_time: 69.384ms







## 关键技术分析

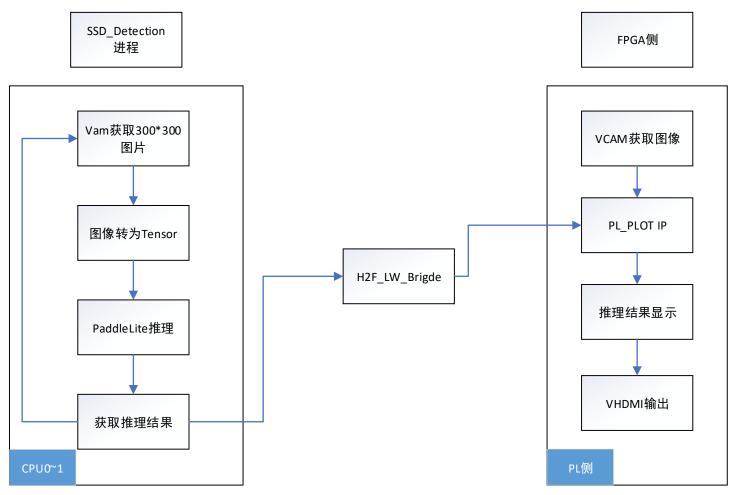
- ·软件程序设计
- ·FPGA工程优化
- -模型优化
- -SDK与PaddLite优化



## 软件程序设计









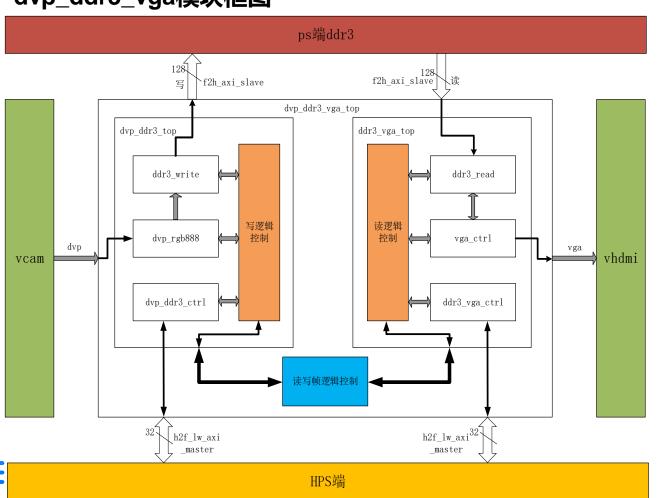
- · 为实现实时推理,我们在分区赛双进程的基础上进一步优化。
- 提出了PL侧绘制预测框的思想,从而可将全部CPU资源用于推理,进一步提升推理刷新速率

## FPGA工程优化-重写dvp\_ddr3以及ddr3\_vga





#### dvp\_ddr3\_vga模块框图



- •重写官方ip,将其整合为dvp\_ddr\_vga 一个ip专门实现**HDMI视频流显示**。
- ·对于该模块,其输入为vcam传来的dvp时序,输出为传输给vhdmi的vga时序。
- •HPS端通过h2f\_lw桥读取dvp\_ddr\_vga模块的**状态寄存器**确认工作状态。

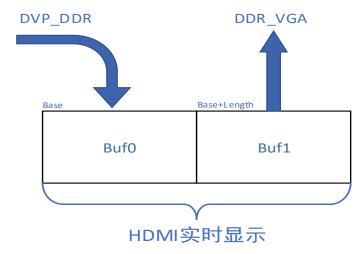
•可自定义添加图片预处理、Resize、PL 端预测框绘制等模块

## FPGA工程优化-双buffer连续写





#### 双buffer连续读写





- •官方的dvp\_ddr3模块HPS请求一帧经过 66ms才能得到数据。
- •我们自己设计的dvp\_ddr3\_vga模块,其无需hps请求控制即可完成对HPS端DDR两个buffer的连续读写,从而无需请求等待的66ms。
- •Vcam的写入与Vhdmi的读取共用双Buf,通过DVP\_DDR以及DDR\_VGA模块互斥读取DDR进行实时视频推流,无需HPS复制。
- •我们**重写了amm驱动**,使其能申请更大的连续内存空间。

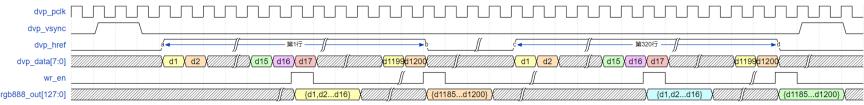


## FPGA工程优化-单通道传输优化

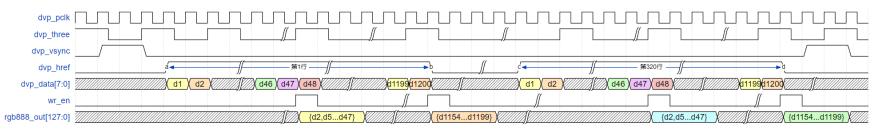




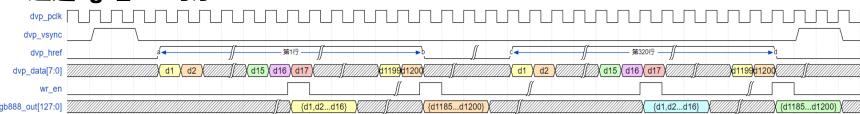




#### 单通道dvp\_rgb888时序



#### 三通道vga\_ctrl时序



•由于测试数据集图片都是灰度图像,三个通道是完全一致的,因此可以训练一个单通道模型、可有效降低模型大小并提升推理速度。

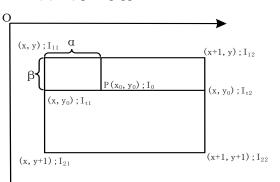
•为了配合HPS端进行单通道模型推理,FPGA端写入HPS的ddr时只写入单通道400\*320大小图片,与原三通道400\*320\*3大小图片相比传输数据量降低两倍。

## FPGA工程优化-重构Resize

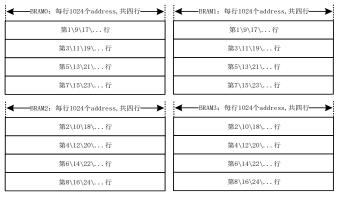




#### 双线性插值示意图

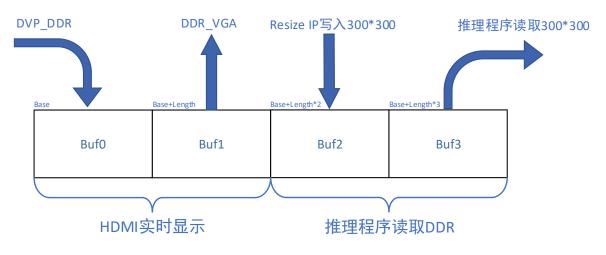


#### FPGA实现



- · 我们在PL<mark>侧</mark>重构了**双线性插值**算法,进行图像的Resize 加速。
- 在FPGA中我们使用四个双口BRAM存储行,分别输出四个像素进行处理。
- 由于在推理刷新过程中,将图片Resize为300\*300大小需要14ms。因此减少resize时间是提高整体推理刷新速率的重要手段。

#### 使用resize模块后的六buffer结构



• 使用Resize的buffer结构如左图所示,其中Buf0~Buf1 存储400\*320的单通道图像,进行视频推流工作。 Buf2~Buf3为Resize后的300\*300单通道图像,使用 FPGA加速后总的图片预处理时间从29ms减少到了 2ms。

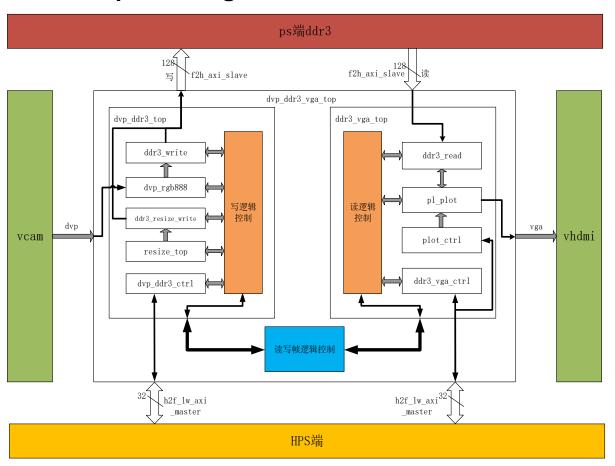
<u>`</u>	Msgs																	
		(Rec	ion: s	m:/tb_r	esize	)												
🔷 dk_in1	1h1										$\neg$		ⅎ					
dvp_vsync	1'h0																	
dvp_href	1'h0										$\neg$							
⊢🔷 dvp_data	8'he0	12		13			14			(15			(16			17	X	18
⊢🔷 img_hs_cnt	11'd0	162		163			(16	1		(165			166			167	Ι	16
img_vs_cnt	11'd14	319																
	1'h1	$\Box$	л	╧	┚				┚		┚		л.	╧	J.			
	1'h0		_										_					
post_img_href	1'h0														щ			
	8'h66	6e	(6f	70	72	73	74	76	77	78	7a	(7b	(7c	) 7e		70		
	11'd0	287	288	289	29	291	29	2 293	294	295	296	297	298	299	300	(0		
<b>├</b> - <b>◇</b> post_img_vs_cnt	11'd12	299														(0		

## FPGA工程优化-PL端叠加推理结果

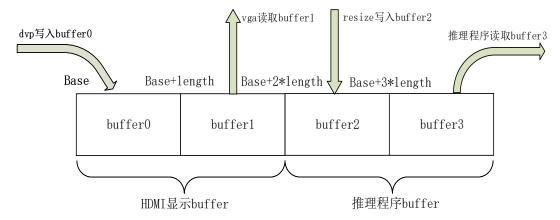




#### 最终的dvp\_ddr3\_vga模块框图



#### 使用resize以及plot模块后的四buffer结构



- 我们在PL侧实现了推理结果的绘制,因此可以将HDMI刷新程序完全由PL端实现,不需要PS端将推理结果叠加到帧图像上,实现HDMI刷新与推理刷新的分离。
- 推理结果的显示由推理程序通过lw总线写入到PL\_Plot的 状态寄存器,完成预测框的实时显示功能。

pre_vga_de	1'h0																																	
<b>-</b> → pix_x_cnt	10'd0	255	256	257 (2	58 2	59 (260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287 2	288
<b>-</b>	10'd 10	10																																
<b>■</b> - <b>♦</b> pix_x	10'd0	255	256	257 (2	58 2	59 [260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287 2	288
<b>+</b> - <b>♦</b> pix_y	10'd 10	10																																
<b></b>	-10'd1	235	236	237 (2	38 2	39 (240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	-1												ij
+	-10'd1	0																				-1												ij

### FPGA工程优化-CNN频率提升

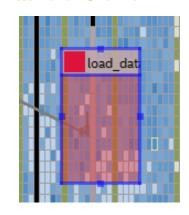




- 移除组件——移除JTAG电路以及无用电路
- 时序约束——使用175Mhz的时序约束文件进行约束
- 时序分析——找到出现不符合时序要求的部分
- 逻辑固化——使用Chip Planner对关键路径逻辑固化
- 时序一致——使f2h\_sdram总线与CNN时钟一致
- 增量编译——优化PL侧模块过多引起的时序不收敛

create\_generated\_clock -name {u0|pll\_1\_cnn|altera\_pll\_i|general[0].gpll~FRACTIONAL\_PLL|vcoph[0]}
-source [get\_pins {u0|pll\_1\_cnn|altera\_pll\_i|general[0].gpll~FRACTIONAL\_PLL|refclkin}] -duty\_cycle
50.000 -multiply\_by 7 -master\_clock {FPGA\_CKLY\_50} [get\_pins {
u0|pll\_1\_cnn|altera\_pll\_i|general[0].gpll~FRACTIONAL\_PLL|vcoph[0] }]
create\_generated\_clock -name {u0|pll\_1\_cnn|altera\_pll\_i|general[0].gpll~PLL\_OUTPUT\_COUNTER|vco0ph[0]}}
-duty\_cycle 50.000 -multiply\_by 1 -divide\_by 2 -master\_clock
{u0|pll\_1\_cnn|altera\_pll\_i|general[0].gpll~PLL\_OUTPUT\_COUNTER|vco0ph[0]}} [get\_pins {
u0|pll\_1\_cnn|altera\_pll\_i|general[0].gpll~PLL\_OUTPUT\_COUNTER|divclk}} ]
create\_generated\_clock -name
{u0|pll\_0\_sdram|altera\_pll\_i|general[0].gpll~PLL\_OUTPUT\_COUNTER|divclk}} -source [get\_pins {
u0|pll\_0\_sdram|altera\_pll\_i|general[0].gpll~PLL\_OUTPUT\_COUNTER|divclk}} -duty\_cycle 50.000
-multiply\_by 1 -divide\_by 3 -master\_clock
{u0|pll\_0\_sdram|altera\_pll\_i|general[0].gpll~FRACTIONAL\_PLL|vcoph[0]} [get\_pins {
u0|pll\_0\_sdram|altera\_pll\_i|general[0].gpll~FRACTIONAL\_PLL|vcoph[0]} [get\_pins {
u0|pll\_0\_sdram|altera\_pll\_i|general[0].gpll~FRACTIONAL\_PLL|vcoph[0]} [get\_pins {
u0|pll\_0\_sdram|altera\_pll\_i|general[0].gpll~FRACTIONAL\_PLL|vcoph[0]} [get\_pins {
u0|pll\_0\_sdram|altera\_pll\_i|general[0].gpll~PLL\_OUTPUT\_COUNTER|divclk}}

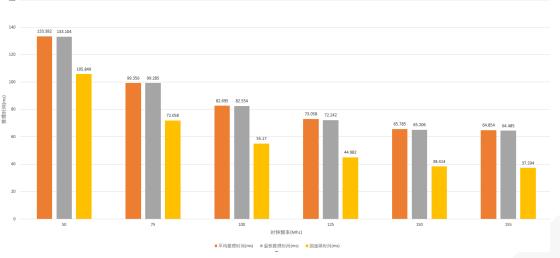
Partition Name	Netlist Type	Color
∨ 看 тор	Source File	
cnn_top:cnn_top_0	Post-Fit	
dvp_ddr3_top_me:dvp_ddr3_top_0	Post-Fit	
receive_top:vhdmi_0	Post-Synthesis	
vcam top:vcam 0	Post-Synthesis	



#### 在82%剪枝率,其他模块最优化条件下

频率(MHZ)	平均推理时间(MS)	最快推理时间(MS)	加速器运行时间(MS)
50	133.382	133.104	105.849
75	99.356	99.285	72.058
100	82.695	82.554	55.170
125	73.058	72.242	44.982
150	65.785	65.206	38.414
155	64.854	64.485	37.294





## 模型优化-敏感度分析剪枝

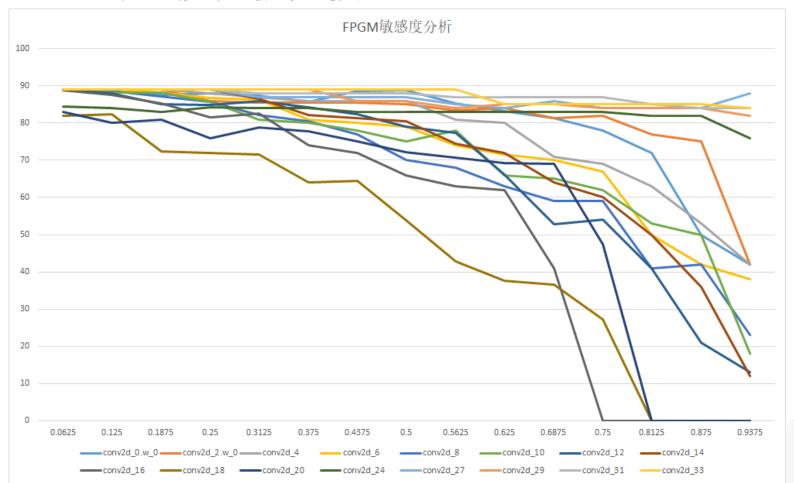




## 敏感度分析剪枝

在剪枝率步进的选择上,由于加速器是**8通道**的,因此我们选择**8**为步进进行敏感度分析(使得通道数为8的整数倍)。

注意:在通道数较少时,对于剪枝率选择应满足剩余通道数是8的整数倍,避免多余通道使加速器多运算一次。





## 模型优化-基于L1和FPGM法则的综合剪枝策略





## 基于L1和FPGM法则的综合剪枝策略

在参考相关综合剪枝论文的基础上,提出在计算**卷积层通道重要性**时选择 FPGM和L1重要性评估方法加权得到重要性分数。

在训练集与测试集划分比例为9:1情况下比较剪枝策略优化前后的mAP指标, 发现mAP由原有的86%提高到91%, 证明混合不同剪枝策略能够一定程度上 在相同压缩率下达到更高的检测精度。

```
reduce dims = [i for i in range(len(value.shape)) if i != pruned axis]
l1norm = np.mean(np.abs(value), axis=tuple(reduce dims))
if groups > 1:
   l1norm = l1norm.reshape([groups, -1])
   l1norm = np.mean(l1norm, axis=1)
# sorted idx = l1norm.argsort()
dist sum list = []
for out i in range(value.shape[0]):
   dist sum = self.get distance sum(value, out i)
   dist sum list.append(dist sum)
scores = np.array(dist sum list)
if groups > 1:
   scores = scores.reshape([groups, -1])
   scores = np.mean(scores, axis=1)
sorted idx = (scores/2+l1norm/2).argsort()
```



## 模型优化-输出层缩放与单通道模型训练





## 输出层缩放

修改mobilenet\_v1.py的源码对6个输出 卷积层的输出通道数进行缩放,从而达到 输出层剪枝的目的。

```
output_prune_ratios=np.array([0.5,0.5,0.5,0.5,0.5,0.5]) dtype=float)
```

```
elif(i==4):
       tmp = self.add_sublayer(
        "conv5_" + str(i + 1),
           sublayer=DepthwiseSeparable(
               in_channels=int(512*prune_ratios[10]),
               # in channels=128,#128
               out_channels1=int(512*prune_ratios[10]),
              out channels2=int(512*output prune ratios[0]),
               num_groups=int(512*prune_ratios[10]),
               stride=1.
               scale=scale,
               conv_lr=conv_learning_rate,
               conv_decay=conv_decay,
               norm_decay=norm_decay,
               name="conv5_" + str(i + 1)))
       self._update_out_channels(int(512*output_prune_ratios[0]), len(self.dwsl), feature_maps)
sublayer=DepthwiseSeparable(
   in channels=int(512* output prune ratios[0]* scale),
   out_channels1=int(512*output_prune_ratios[0])
   out_channels2=1024*prune_ratios[11],
   num groups=int(512*output prune ratios[0]).
```

## 单通道模型

- 1.对第一个输入卷积层的通道数进行修改
- 2.对图像读取的算子进行修改

(**Decode**、Resize、NormalizeImage)

```
self.conv1 = ConvBNLayer(
    in_channels=1,

im_copy = cv2.cvtColor(im, cv2.COLOR_BGR2GRAY)

im= np.zeros((im.shape[0], im.shape[1], 1), np.uint8)

im[:,:,0]=im_copy

Decode:
    像素值单通道填充
    产生输入数组
```



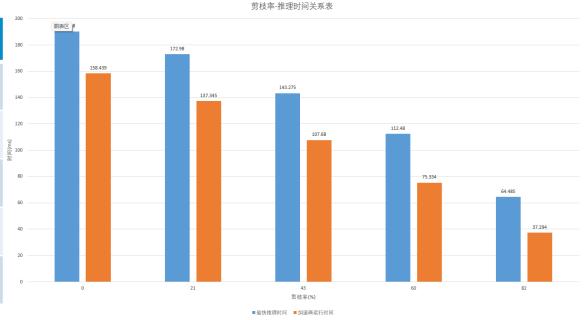
## 模型优化-模型剪枝率提升效果图





#### 我们最终实现了82%的剪枝率,提升效果如下所示

剪枝率(%)	输出层缩放 全量化	GFLOPs	最快推理时间(MS)	加速器运行时间(MS)
0	是	5.09	190.328	158.439
21	是	3.99	172.98	137.345
43	是	2.88	143.275	107.680
60	是	1.97	112.48	75.334
82	是	0.93	64.102	37.294





## SDK优化-数据重排优化





### 1.开启O3优化

```
set(CMAKE_CXX_FLAGS "-march=armv7-a -mfloat-abi=hard -mfpu=neon ${CMAKE_CXX_FLAGS} -03")
set(CMAKE_C_FLAGS "-march=armv7-a -mfloat-abi=hard -mfpu=neon ${CMAKE_C_FLAGS} -03")
```

### 2.去除无用memset()操作

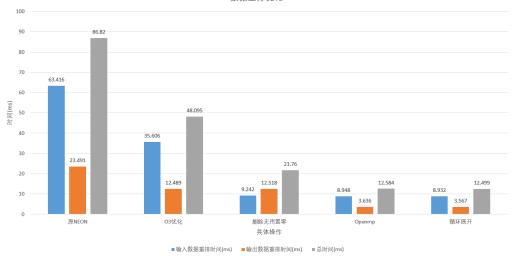
#### 3.使用OPENMP优化

#### 4.循环展开

```
#pragma unroll(32)
for (int i = 0; i < up_round(c, INPUT_EXTEND_SCALE); i++)
{
   tran_8((uint8_t *)din + i * area * INPUT_EXTEND_SCALE, (uint8_t *)dout + i * area * INPUT_
}</pre>
```

具体操作	输入数据重排时 间(MS)	输出数据重排时 间(MS)	总时间(MS)
原NEON	63.416	23.491	86.82
O3优化	35.606	12.489	48.095
删除无用 置零	9.242	12.518	21.76
OPENMP	8.948	3.636	12.584
循环展开	8.932	3.567	12.499

#### 数据重排优化



## PaddLite优化-修复PaddLite能耗模式BUG





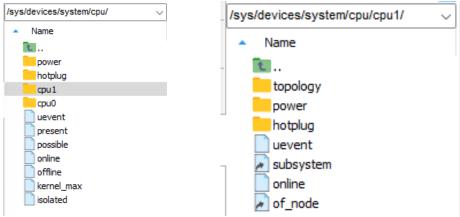
在模型加载时会出现如下报错:查询CPU0状态失败,切换到NO BIND模式

```
[I 5/24 12:50:20.521 ...nna/Paddle-Lite/lite/core/device_info.cc:509 check_cpu_online] Failed to query the online statu
e of CPU id:0
[I 5/24 12:50:20.521 ...nna/Paddle-Lite/lite/core/device_info.cc:514 check_cpu_online] CPU id:0 is offine
[W 5/24 12:50:20.522 ...nna/Paddle-Lite/lite/core/device_info.cc:1352 SetRunMode] Some cores are offline, switch to NO
BIND MODE
```

经过调试,定位PaddleLite绑定核心源码位于 device info.cc文件下

```
ool bind threads(const std::vector<int> cpu ids) {
 int thread num = cpu ids.size();
omp set num threads(thread num);
 std::vector<int> ssarets(thread num, 0);
pragma omp parallel for
 for (int i = 0; i < thread num; i++) {
  ssarets[i] = set sched affinity(cpu ids);
 for (int i = 0; i < thread num; <math>i++) {
  if (ssarets[i] != 0) {
    LOG(ERROR) << "Set cpu affinity failed, core id: " << cpu ids[i];
#else // ARM WITH OMP
 std::vector<int> first cpu id;
 first cpu id.push back(cpu ids[0]);
 int ssaret = set sched affinity(first cpu id);
  LOG(ERROR) << "Set cpu affinity failed, core id: " << cpu ids[0];
endif // ARM WITH OMP
```

经过分析,PaddleLite通过查询/sys/devices/system/cpu/cpu0目录下的online文件,来确认核心是否可用并绑定。但由于HPS中CPU0为主CPU,默认在线,文件中并未含有online文件。以至于PaddleLite误以为CPU0不可用,将进程绑定在CPU1上,从而限制系统性能



通过修改源码,使其绑定双核心并利用双线程运行,充分利用硬件资源,修改前后推理刷新时间提高了10ms。

```
input_organize_time:10.116718ms
fpga_time:37.284771ms
output_organize_time:7.526860ms
------graph end-----
FPS_run:87.175ms
```



input\_organize\_time:10.141111ms
fpga\_time:37.325283ms
output\_organize\_time:3.697364ms
-----graph end----FPS run:77.784ms



### PaddLite优化-预处理优化





主办方提供数据集为**灰度图**, RGB888格式, 且PaddleLite官方Demo仅支持三通道输入, 将框架输入改为**单通道**有如下好处:

- 1. dvp\_ddr IP输入为单通道,其相邻像素间地址连续,HPS从DDR读取原始图像数据将节省读取时间。
- 2. 图片需要Resize为300\*300送入模型推理, 单通道图像可减少Resize时间。
- 3. 图像数据需转换为(-1,1)范围内的Tensor格式,即uint8->float32,单通道将减少转换时间。
- 4. NHWC->NCHW通道转换, PaddleLite 的输入特征图为NCHW格式,需进行格式转换以及归一化操作,单通道可以缩短通道转换以及归一化时间。

我们通过修改单通道预处理过程,并且 FPGA加速Resize,将数据预处理时间从 29ms缩短至2ms,如下图所示。

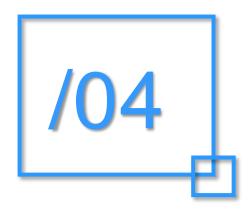
预处理过程	修改前时间(MS)	修改后时间(MS)
RESIZE单通道优化	16.284	5.097
UINT8->FLOAT32	4.326	0.882
NHWC->NCHW	7.501	1.289
RESIZE PL加速	5.097	0
预处理总时间	28.977	2.171











## 总结与补充







## PL侧相关优化

- 重写DVP\_DDR, DDR\_VGA IP, 修改其 为单通道输入输出,减少数据传输 时间。
- 重写架构,视频推理完全由PL控制, 不受HPS性能影响。
- 使用PL对推理过程中的Resize进行加速,完成双线性插值算法的FPGA实现。
- 使用PL设计了PL\_Plot IP,利用 FPGA完成预测框的绘制,从而可将 全部CPU资源用于推理,提升推理刷 新速度。
- 使用增量编译、逻辑固化、时序约束、去除无用组件如JTAG等。

## 模型训练优化

- 剪枝与量化结合训练,加快模型推理速度
- 基于敏感度的剪枝分析,达到81.6%的剪枝率
- 卷积层缩放,进一步减小模型大小
- 卷积连续,算子融合,能够利用CNN 加速器加速卷积运算
- 数据增强:增强轮廓特征、水平变换, 提高模型识别准确度
- 单通道模型修改,为预处理时间优化铺垫

## PS侧相关优化

- 修改PaddleLite架构,使其完成单 通道推理,将预处理时间缩短至2ms。
- 修改PaddleLite底层源码,使其能够利用双核心推理,提升推理刷新速度。
- 利用NEON指令集完成数据重排加速, 从83ms缩短至13ms。
- 修改Cmake文件,开启03优化,使用OpenMP优化进程。
- 利用Linux的进程调度功能,使用 CPU核心绑定,进程优先级调整,利 用脚本文件提升系统推理性能。



## 补充

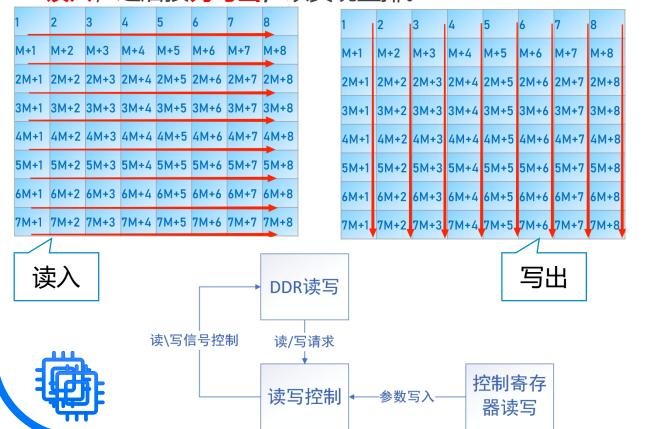
## PL端数据重排





## verilog实现

我们的基本思路是在读写控制中设法构建一个**8\*8的数据矩阵**,对于该矩阵,数据**按行 读入**,之后按**列写出,**以实现重排。



## OpenCL实现

OpenCL的实现主要起验证作用,由于经过分析,我们认为PL端数据重排并不能比ARM侧使用NEON指令集加速,因此最终并未选用PL数据重排

```
double time1 = getCurrentTimestamp();
status = clEnqueueWriteBuffer(queue, dev_src, CL_TRUE, 0, sizeof(cl_int)*src_size, src, 0, NULL, NULL);
status = clEnqueueWriteBuffer(queue, dev_in_c, CL_TRUE, 0, sizeof(cl_int), in_c, 0, NULL, NULL);
status = clEnqueueWriteBuffer(queue, dev_in_h, CL_TRUE, 0, sizeof(cl_int), in_h, 0, NULL, NULL);
status = clEnqueueWriteBuffer(queue, dev_in_w, CL_TRUE, 0, sizeof(cl_int), in_w, 0, NULL, NULL);
```

total\_time is:1.289770 ms root@socfpga:~/zhuanzhi#





# 感谢观看

## 汇报结束 请评委专家提问!

A组 CICC1577 冰糖葫芦儿

