# XJTU-Tripler:
# DNN Accelerator for Autonomous Agents

Boran Zhao, Wenzhe Zhao, Tian Xia, Fei Chen, Long Fan, Pengchen Zong,

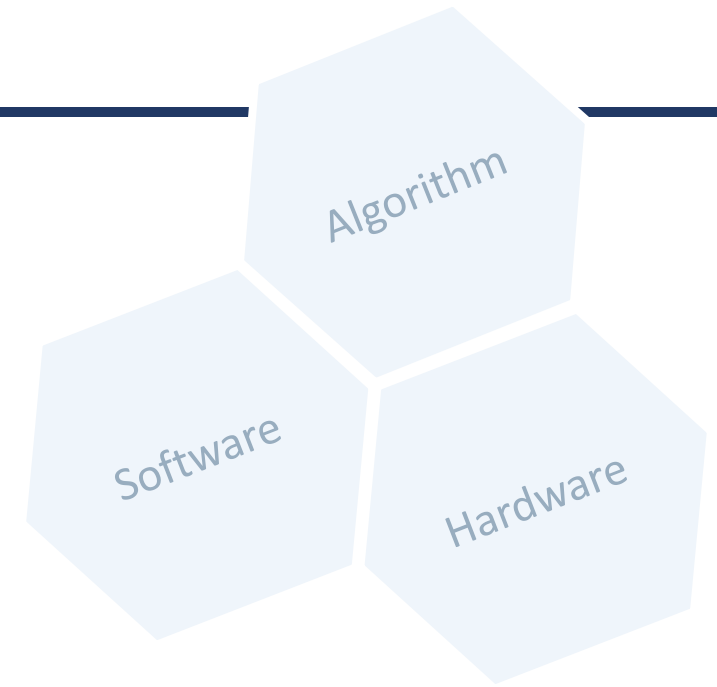Yadong Wei, Zhijun Tu, Zhixu Zhao, Zhiwei Dong, and Pengju Ren*

XJTU-Tripler Group, IAIR, Xi'an Jiaotong University

(pengjuren@xjtu.edu.cn)*

2019-06-05

# Outline

- Background & Motivation

- Network selection and quantization for single object detection

- HiPU: a General-purpose Fixed-point DNN accelerator on Xilinx FPGA

- Data flow of HiPU

- Optimization of HiPU and Zynq system

- Performance analysis

- Roadmap

- Demo 1: Dual-channel real-time object detection on Ultra96(Xilinx Zu3)

- Demo 2: Autonomous agents cooperation on TX2 now, on Ultra96 in near future
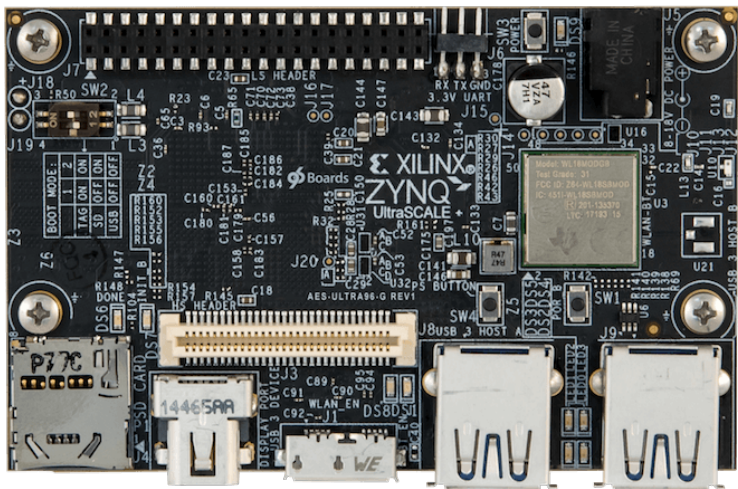
Algorithm

Software

Hardware

# Background & Motivation

## Background

The DAC system design contest aims to implement single-object detection algorithm on various class of pictures, including car, person, bicycle, and other indicated objects.

The best technology supporting such detection task is Deep Neural Network, which requires tremendous amount of computation.

Ultra96(Xilinx ZU3) is an outstanding platform to implement DNN accelerator with low power.



| ZU3 resource | |
|------|------|
| LUT | 70K |
| FF | 141K |
| BRAM | 216 |
| DSP | 360 |

## Motivation

- Algorithm

Customized YOLO is exploited to trade off the accuracy and energy.

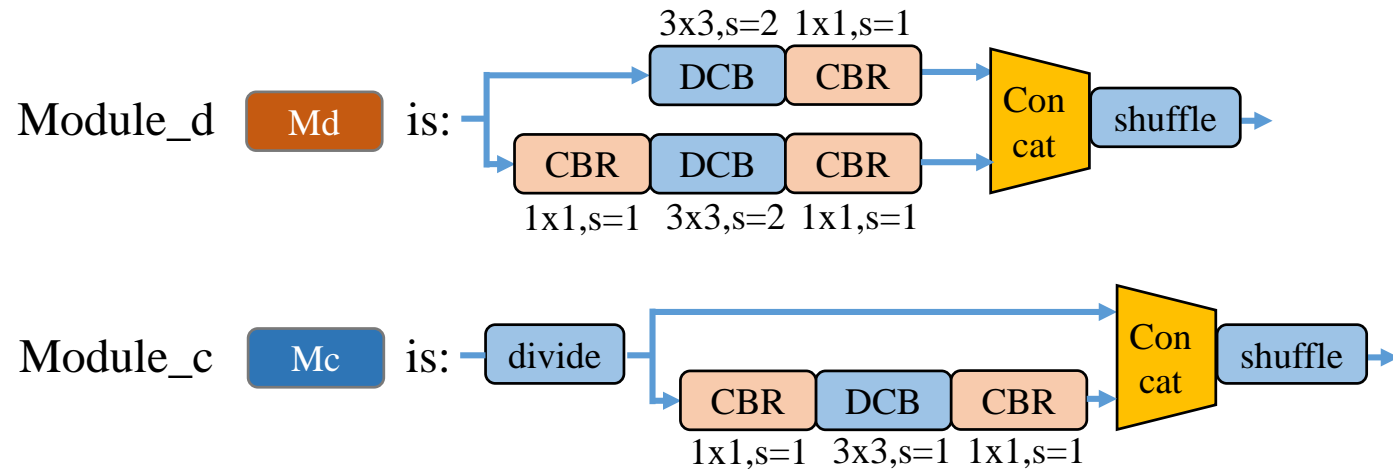| 1 | Replace the feature extraction with ShuffleNet V2 |
|---|---|
| 2 | Exploit 8-bit quantization and fine tune technique |

- FPGA implementation

A general-purpose DNN accelerator, HiPU, is implemented to support all of the special network operators.

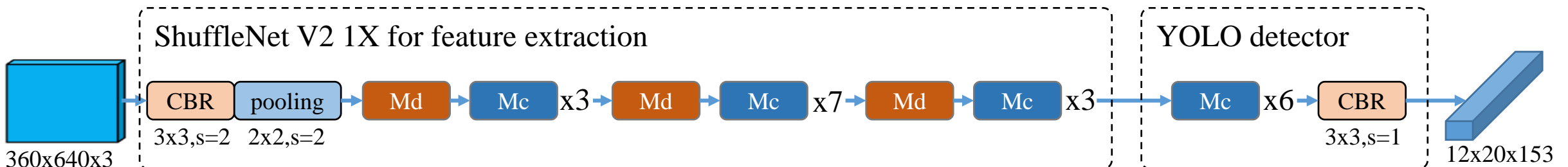| 1 | HiPU supports CONV, depth-wise CONV, padding, pooling, RELU, channel shuffle/concat/divide |
|---|---|
| 2 | HiPU supports inter-layer optimization to reduce DDR bandwidth by the embedded RISC-V core |
| 3 | The system concurrently processing LD/CONV, PS/PL |

# Network selection for single object detection: YOLO with ShuffleNet V2

At first, we define **CBR** as the combination of **Conv** **BN** **ReLU** , and **DCB** is the combination of **Depthwise conv** **BN** .

Then, 2 typical bottlenecks of ShuffleNet V2 can be described as follows:

Module_d **Md** is:

Module_c **Mc** is:

Our final detection network is:

In order to reduce the parameter quantity of conventional YOLO, ShuffleNet V2 is exploited to replace the feature extraction layers.
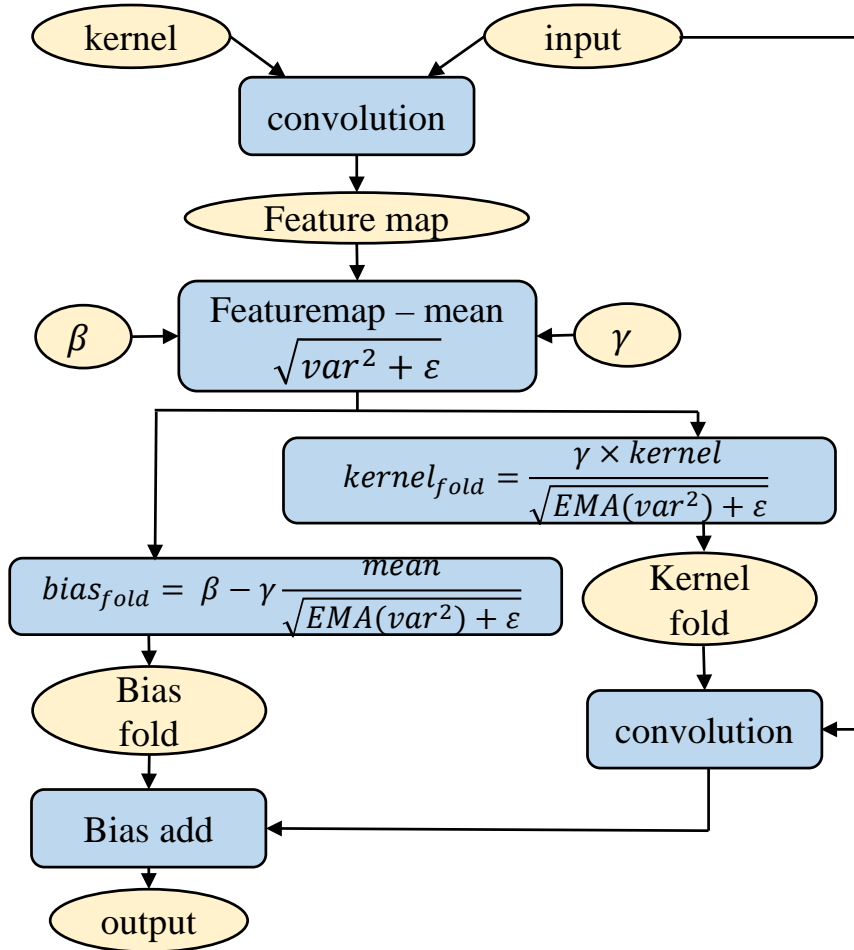
At first, the parameters of feature extraction layers are pre-trained on ImageNet classification dataset. Parameters before the first Module_c was reloaded and trained when it was fine-tuned on DAC dataset.

At last, YOLO detector on DAC19 dataset was trained.

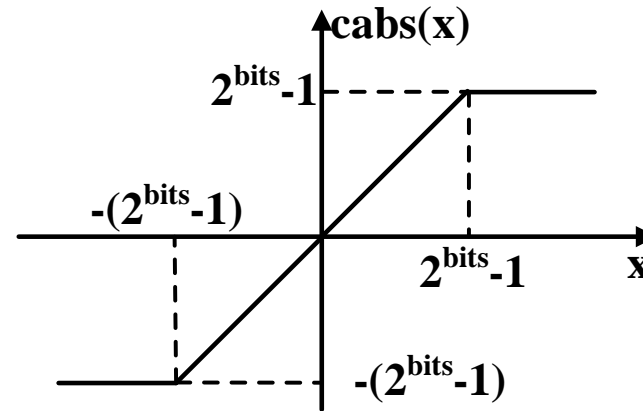To improve the detection accuracy, the input image remains original size.

# Quantization of Neural Network

**Before quantization, make fusion of BN layer:**

kernel → convolution ← input

convolution → Feature map

Feature map → $\dfrac{Featuremap - mean}{\sqrt{var^2 + \varepsilon}}$ (with $\beta$ and $\gamma$ inputs)

$$kernel_{fold} = \frac{\gamma \times kernel}{\sqrt{EMA(var^2) + \varepsilon}}$$

$$bias_{fold} = \beta - \gamma \frac{mean}{\sqrt{EMA(var^2) + \varepsilon}}$$

Kernel fold → convolution

Bias fold

Bias add → output

Then, $output = input \times kernel_{fold} + bias_{fold}$

---

**8bit symmetry quantization:**

$$ratio_{index} = \left\lceil \frac{2^{bits} - 1}{max(tensor)} \right\rceil \quad (1)$$

$$ratio = 2^{ratio\_index} \quad (2)$$

$$quan(x) = [x] \quad (3)$$

$$cabs(x) = \begin{cases} 2^{bits} - 1 & x > 2^{bits} - 1 \\ -(2^{bits} - 1) & x < -(2^{bits} - 1) \\ x & else \end{cases} \quad (4)$$

$$output = \frac{quan(cabs(ratio_a \times input)) \times quan(cabs(ratio\_k \times kernel))}{ratio\_a \times ratio\_k} + \frac{quan(cabs(ratio\_b \times bias))}{ratio\_b}$$

---

After quantization, the network should be fine-tuned to restore accuracy.

Once the fixpoint parameter is determined, transform the parameters to FPGA format.

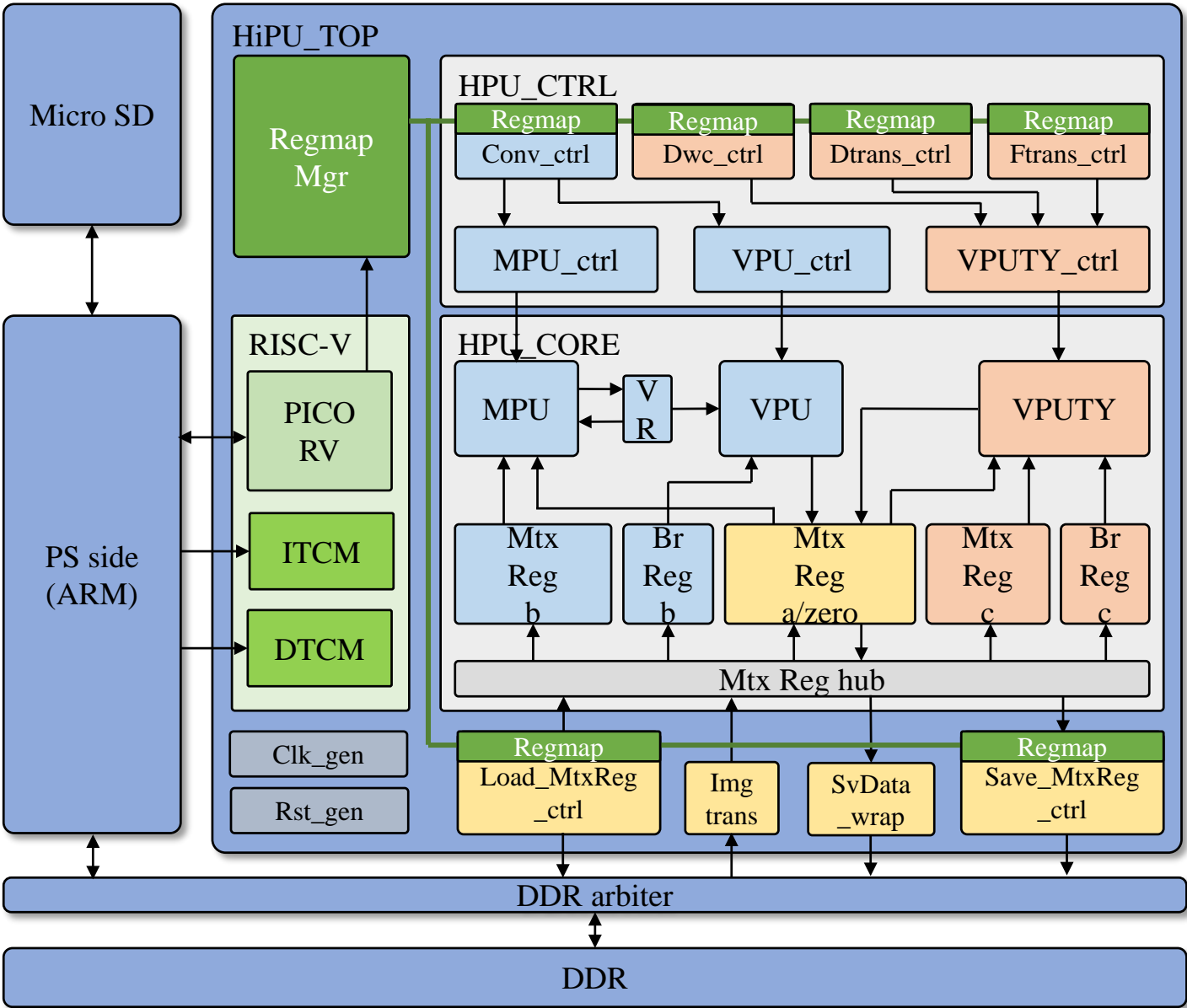The size of our proposed network is: **Weight: 1.94MB, Bias: 78KB**

# HiPU: a General-purpose Fixed-point DNN accelerator on Xilinx FPGA

HiPU is a general purpose processing unit, which can

1. calculate Matrix mul-add by MPU;

2. calculate Vector mul/add/cmp/shift operation by VPU or VPUTY;

3. calculate Scalar operations and controlling logic by embedded RISC-V core.

DNN operator can be processed by abovementioned modules:

| DNN operator | Module of HiPU |
|---|---|
| Convolution | MPU, VPU |
| Depth-wise CONV | VPUTY |
| Padding | VPUTY, RISC-V schedule |
| Pooling | VPUTY |
| RELU | VPU, VPUTY |
| Channel shuffle | VPU, VPUTY |
| Concat/divide | RISC-V schedule |

The weight, bias, and input image/feature map are stored in DDR.

Since HiPU only processes the data on MtxReg, we should transform data from DDR to MtxReg before calculation, and transform data from MtxReg to DDR after calculation is finished.
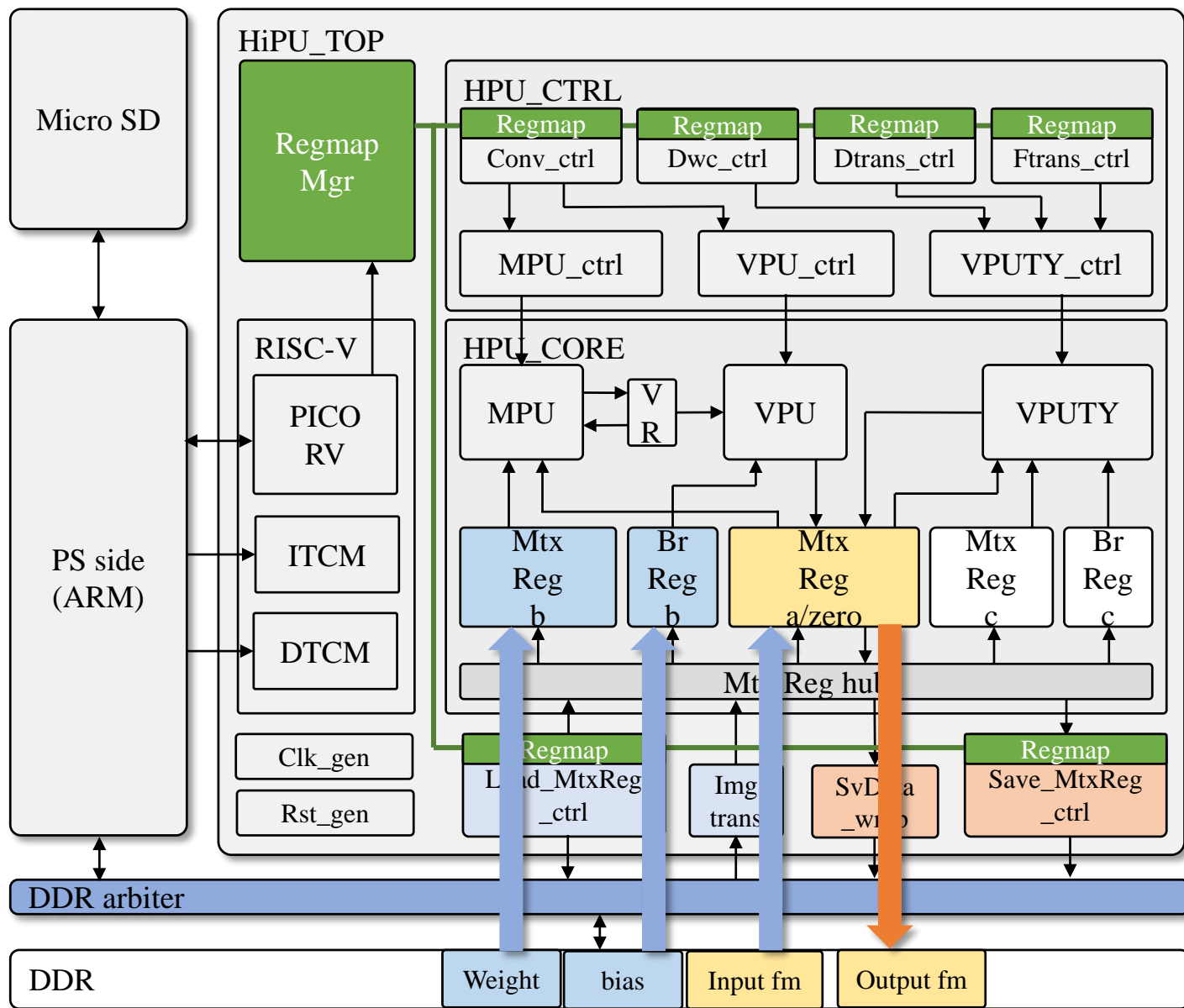
**LDMR:**

The conventional convolution data preparation is shown as the BLUE arrow of right figure.

According to the regmap configuration, LDMR_ctrl module transforms weight, bias, and input fm from DDR to MtxReg.

**SVMR:**

The conventional convolution data output is shown as the ORANGE arrow of right figure.

According to the regmap configuration, SVMR_ctrl module transforms output fm from MtxReg to DDR.

# Data flow of HiPU: convolution(CONV) and depth-wise convolution(DWConv)

HiPU supports variety of operators, which can be classified in 4 group: Conv, Depth-Wise Calc(DWC), Data-Trans(Dtrans), and Format-Trans(Ftrans).

Take CONV and DWConv as example:

**CONV**:

Phase1: Conv_ctrl module sends a series of instructions to MPU and VPU;

Phase2: MPU reads weight and input fm from MtxReg;

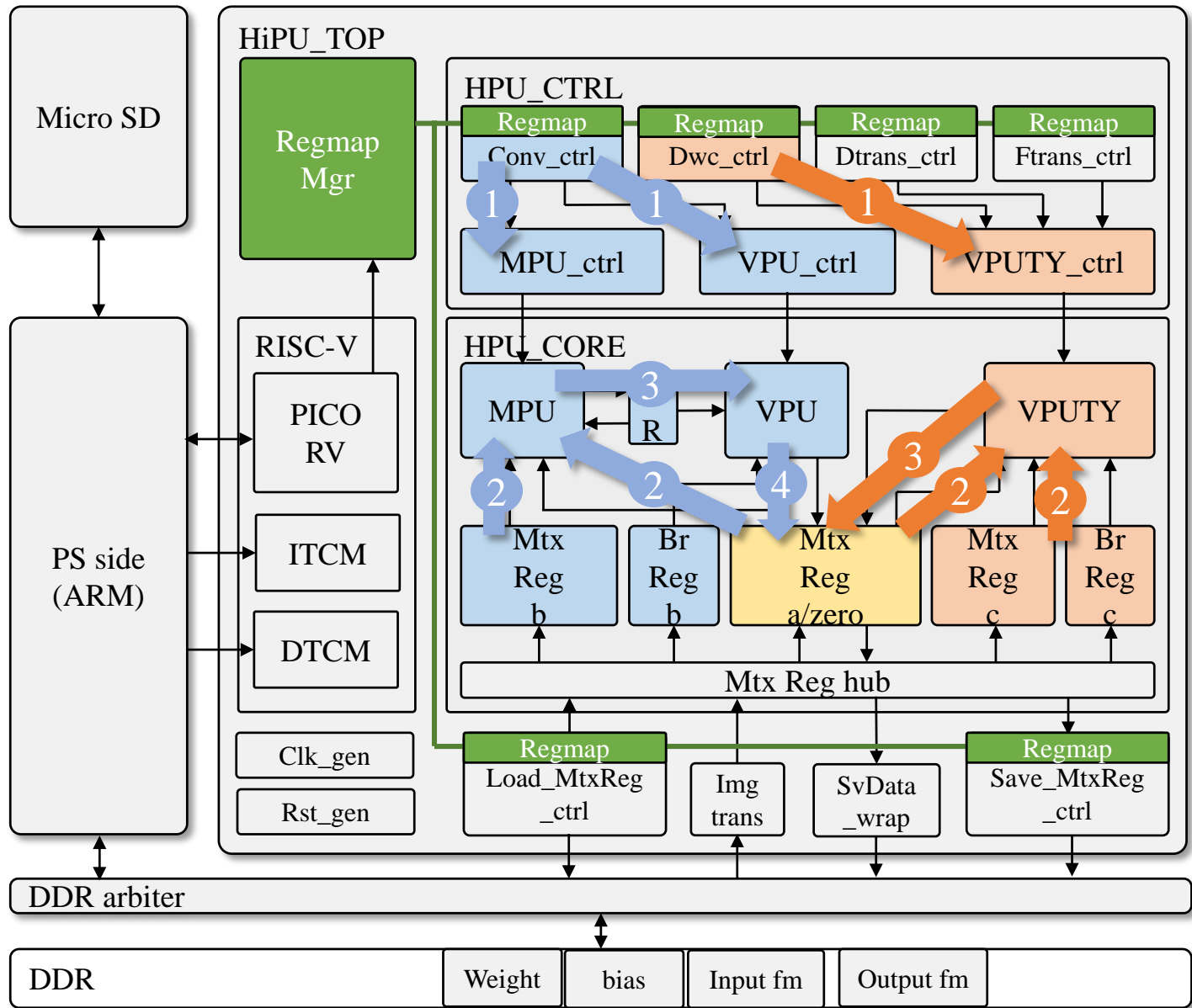Phase3: MPU makes Matix-mul-acc operations and output result to VPU;

Phase4: VPU processes adding bias, ReLU, shuffle operations for every point of output fm, and writes back to MtxReg.

**DWConv**:

Phase1: Dwc_ctrl module sends a series of instructions to VPUTY;
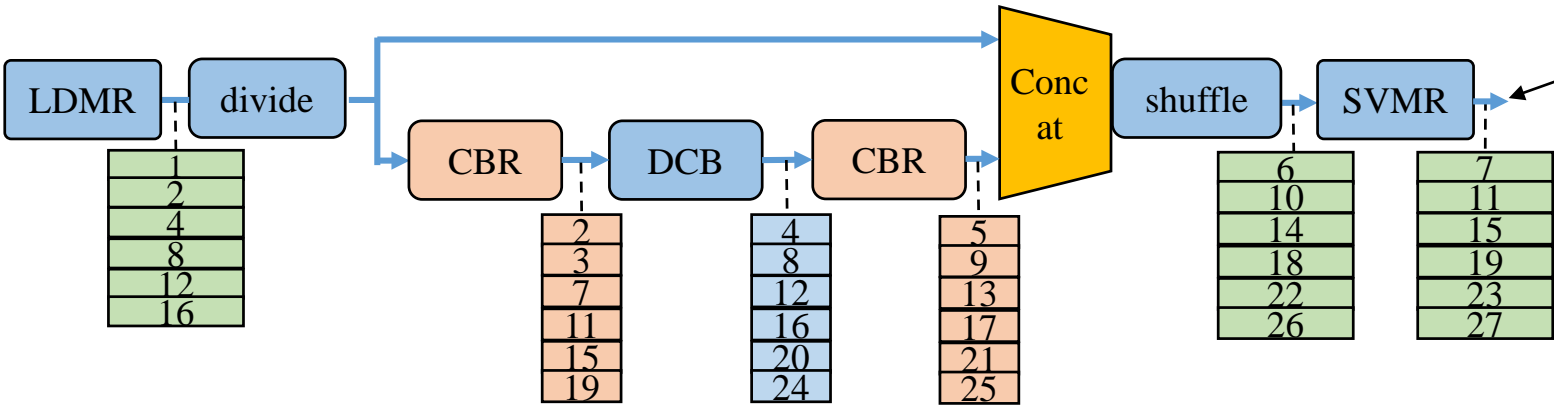
Phase2: VPUTY reads weight and input fm from MtxReg;

Phase3: VPUTY makes Vector-mul-acc, adding bias, ReLU, shuffle operations, and write back to MtxReg.

# Optimizations of HiPU

## Tip 1: Processing all layers of a bottleneck together to reduce DDR bandwidth

Since the fm size at DNN bottleneck is limited, it is reasonable to load fm from a bottleneck starts, and write back fm from the bottleneck ends. A carefully schedule on RISC-V is needed.
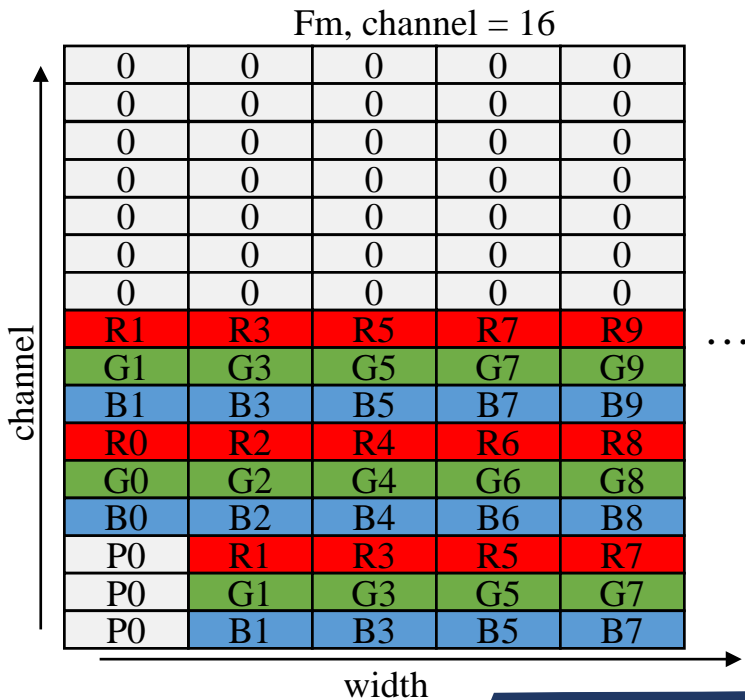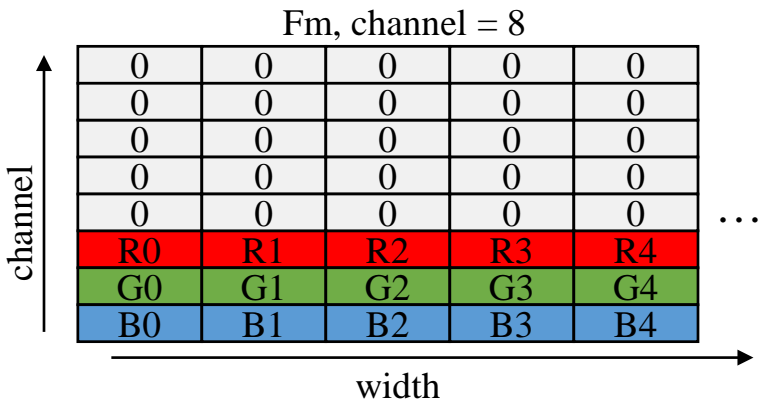
Give module_c as example.

The number of block means the height of fm.

When processing Module_c, DDR read/write operation needs only 1 time.



## Tip 2: Transform input image to increase the HiPU efficiency

Since the channel of input image is only 3, leading to low efficiency of HiPU calculating. HiPU has a special transform module to extend the channel of input image from 3 to 9 if the kernel of Conv1 is 3x3.
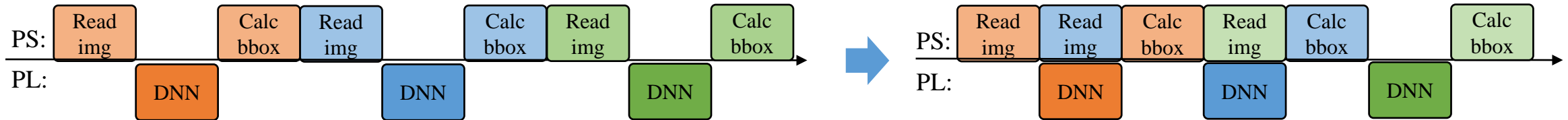
Increase the HiPU efficiency from 0.38 to 0.56 for Conv1

# Optimizations of Zynq system

**Tip 1: Parallel in read image and HiPU processing**

Since HiPU does not have dependence with ARM core on PS, read image can work concurrently with HiPU calculation.



The FPS of processing increases from 30.3Hz to 49.9Hz

**Tip 2: PS computation is accelerated by C functions**

Some time-consuming software routines are accelerated using C code and imported to Python with ctypes:

- Bounding box generation acceleration
  Customized C functions. Precompute the addresses of confidence data and bbox data in the output data of PL. Find the max confidence and the corresponding bbox, then compute the absolute bbox coordinates quickly at run time.

  The time of calc bbox module decreases from 2ms to 0.6ms

- Image reading acceleration
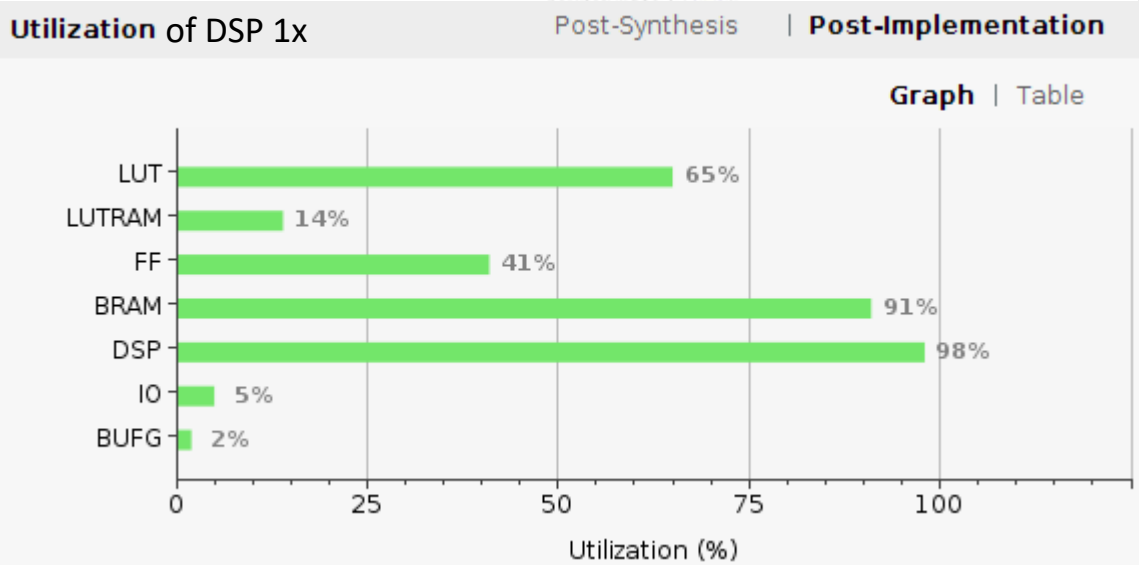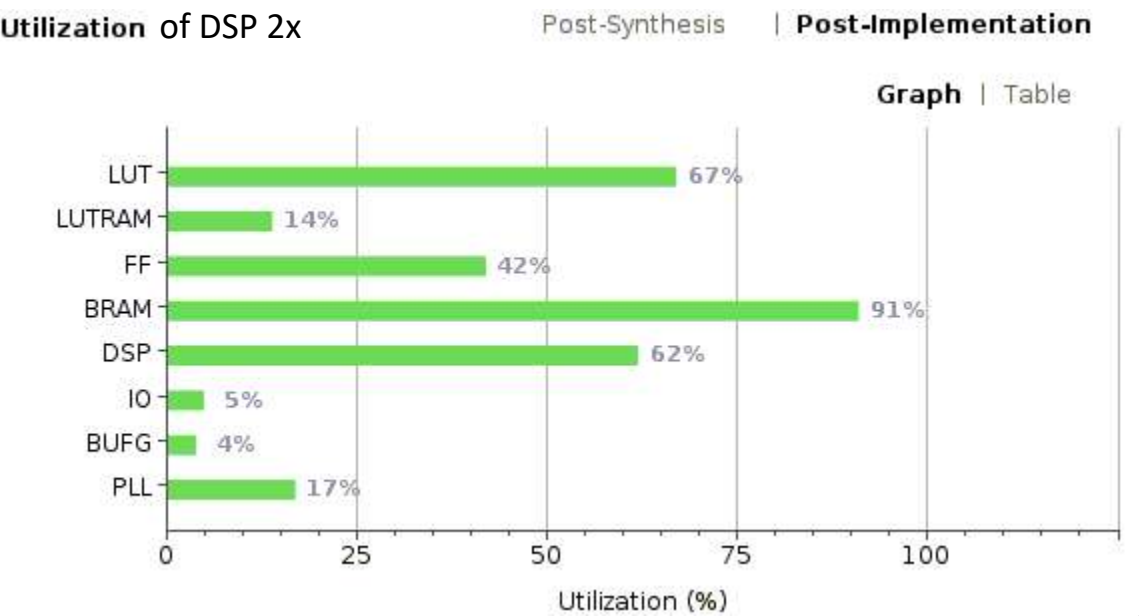  Customized image reading C function using libjpeg.

  This code is just finished, not include in our final version. This should improve the FPS from 49.9Hz to 52Hz.

**Tip 3: Exploit clock gating to reduce power**

In order to reduce the energy consumption of HiPU, close the clock of PL side when HiPU finishes calculation of 1 picture.

# Performance of HiPU

| Characteristics of HiPU | | |
|---|---|---|
| Frequency | | 233MHz |
| Peak Performance | | 268Gops |
| Efficiency | | >80% |
| Developing API | | C/RISC-V ASM |
| Operator support | CONV | Kernel size (1 – 7) |
| | Stride | Stride (1 – 256) |
| | Dilation | Dilation size(0 – 255) |
| | Padding | Padding size(0 – 255) |
| | Pooling | Range: (1-7), Global; Type: Min/Max/Ave |
| | Depth-wise Conv | Kernel size (1-7) |
| | concat/divide | Any |
| | Channel shuffle | Just support Grp=2 |
| | DeConv | Will support in future |
| | Ele-wise Add/Mul | Support |
| | Activation | Only support ReLU |
| Reliance | Could be deployed on any kind of Xilinx FPGA (ZU, KU, VU, 7series) , do not rely on Zynq. | |



Utilization of DSP 2x — Post-Synthesis | Post-Implementation — Graph | Table

Utilization (%): LUT 67%, LUTRAM 14%, FF 42%, BRAM 91%, DSP 62%, IO 5%, BUFG 4%, PLL 17%



Utilization of DSP 1x — Post-Synthesis | Post-Implementation — Graph | Table

Utilization (%): LUT 65%, LUTRAM 14%, FF 41%, BRAM 91%, DSP 98%, IO 5%, BUFG 2%

# Performance analysis

Test Environment： SD-Card type is SanDiskA1-U1-Class10-64GB, 40260 test images in size of 640*360

| System Freq, DSP2x | Latency(ms) | FPS | Power(W) | Energy(KJ) | Efficiency |
|---|---|---|---|---|---|
| 75MHz | 62.89 | 15.90 | 5.88 | 19.43 | 0.37J/pic |
| 166MHz | 21.30 | 46.95 | 8.96 | 10.02 | 0.19J/pic |
| 200MHz | 21.20 | 47.17 | 9.08 | 10.11 | 0.19J/pic |

| System Freq, DSP1x | Latency(ms) | FPS | Power(W) | Energy(KJ) | Efficiency |
|---|---|---|---|---|---|
| 133MHz | 29.80 | 33.56 | 6.59 | 10.21 | 0.2J/pic |
| 150MHz | 26.90 | 37.17 | 6.89 | 9.7 | 0.18J/pic |
| 166MHz | 21.30 | 46.95 | 7.76 | 8.68 | 0.17J/pic |
| 200MHz | 21.30 | 46.95 | 7.79 | 8.71 | 0.17J/pic |
| 233MHz | 21.00 | 47.62 | 7.83 | 8.64 | 0.16J/pic |
| 200MHz(Clk Gating) | 20.84 | 47.98 | 7.21 | <u>7.89</u> | <u>0.15J/pic</u> |
| 200MHz(CG, opti conv1) | 19.93 | 50.15 | 7.40 | **7.75** | **0.15J/pic** |

Notice: The energy data in the table are equivalent to the scale of the competition's test set (52,500 images).
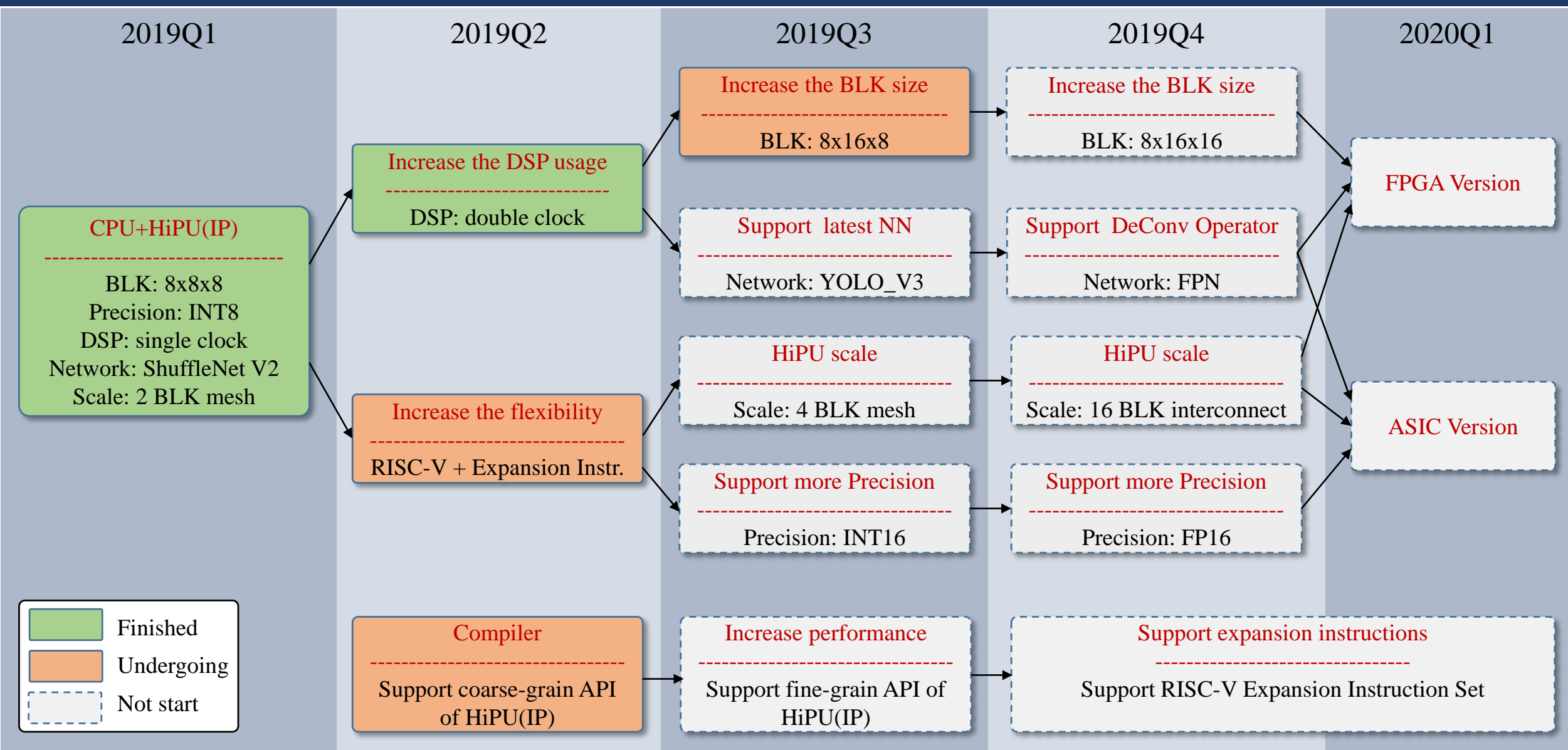
**Conclusion 1:**
When system clock is higher than 166MHz, the FPS is difficult to increase. HiPU is not the bottleneck of the system.
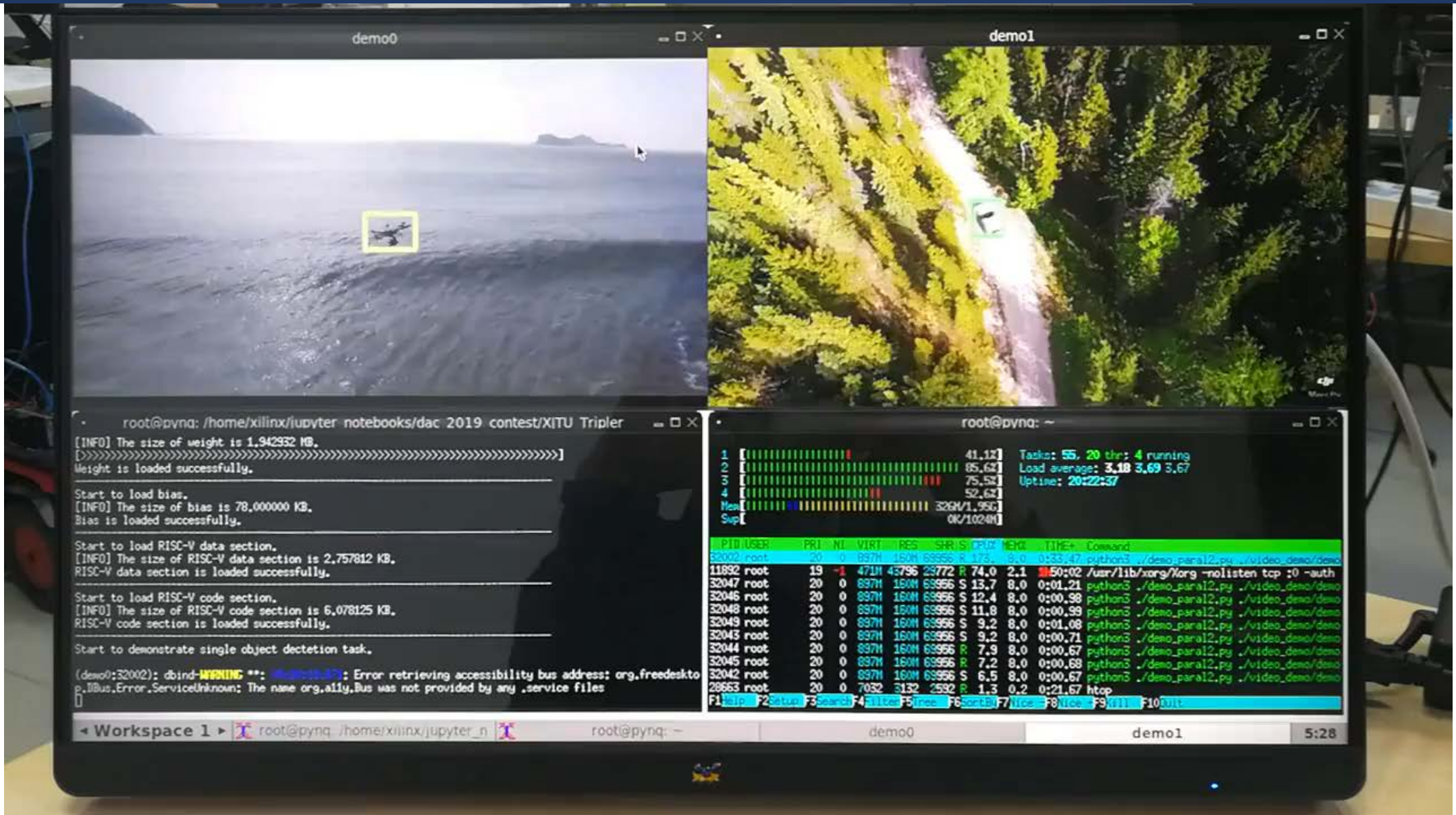
**Conclusion 2:**
We tried to make DSP work on double clock to reduce the FPGA resource. However, the energy is higher than the equivalent single clock DSP version.

According to the performance analysis, we choose to submit the version of DSP1x 200MHz, with clock gating, optimized Conv1 schedule.

# Roadmap

# Demo 1: Dual-channel real-time object detection on Ultra96(Xilinx Zu3)

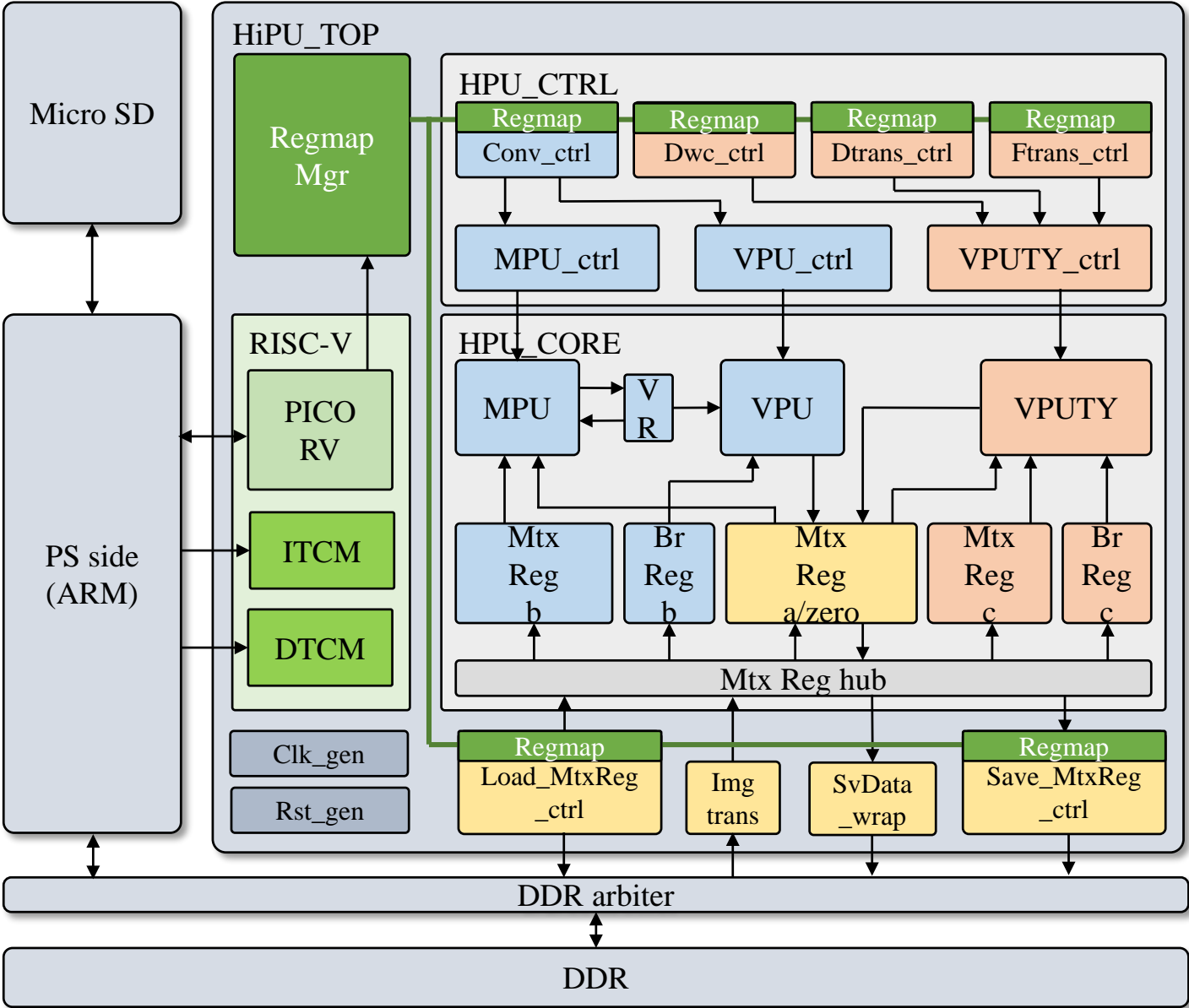# Q & A

Thank you for your attention

# HiPU: a General-purpose Fixed-point DNN accelerator on Xilinx FPGA

| Characteristics of HiPU | | |
|---|---|---|
| Frequency | | 233MHz |
| Peak Performance | | 268Gops |
| Efficiency | | >80% |
| Developing API | | C/RISC-V ASM |
| Operator support | CONV | Kernel size (1 – 7) |
| | Stride | Stride (1 – 256) |
| | Dilation | Dilation size(0 – 255) |
| | Padding | Padding size(0 – 255) |
| | Pooling | Range: (1-7), Global; Type: Min/Max/Ave |
| | Depth-wise Conv | Kernel size (1-7) |
| | concat/divide | Any |
| | Channel shuffle | Just support Grp=2 |
| | DeConv | Will support in future |
| | Ele-wise Add/Mul | Support |
| | Activation | Only support ReLU |
| Reliance | | Could be deployed on any kind of Xilinx FPGA (ZU, KU, VU, 7series) , do not rely on Zynq. |

# Network selection for single object detection: YOLO with ShuffleNet V2