



赛灵思技术日
XILINX TECHNOLOGY DAY

用赛灵思 FPGA 加速机器学习推断

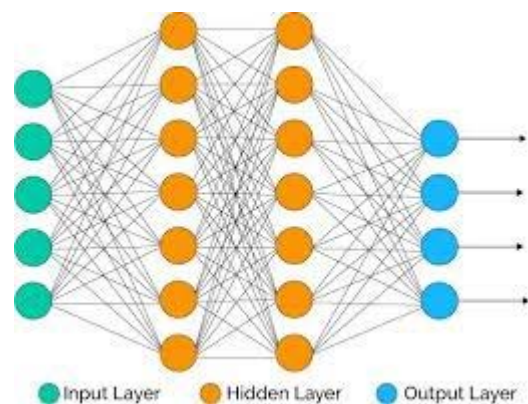
張帆

資深全球AI方案技術專家

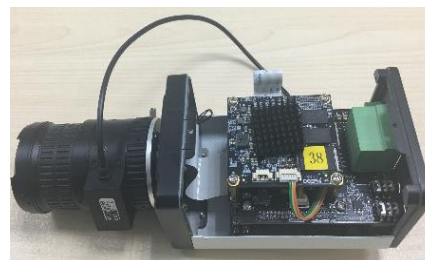
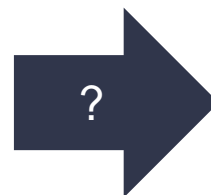
2019年4月25日

 XILINX®

AI在edge端应用中落地面临的挑战



= 43 TOPS



1080p Object Detection (SSD) @ 30 FPS

< 10W, < 50 ms latency, <\$50

Challenges Are Opportunities

<5%

ML expertise
with FPGA

25%

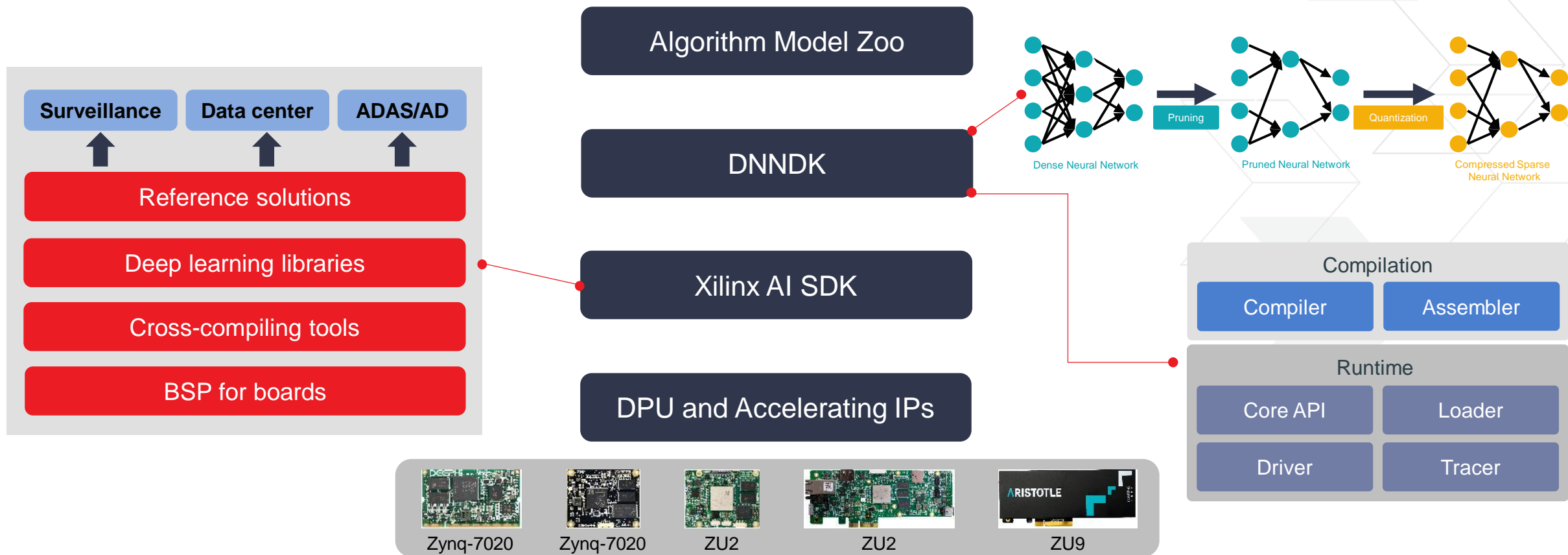
ML expertise
with GPU/CPU/DSP

>70%

No expertise
Looking for near
“turn-key” solutions

Number of Opportunities

Xilinx AI Solution



Algorithm Model Zoo

69 Models in the zoo,
in which 34 models
are already deployed.

Network	Backbone	DPU Deployment	Input Size	OPs	Paras	Training Set	Val Set	Eva Norm Float	Eva Norm Fixed
Resnet50_v1	Resnet	Yes	224*224	7.7G	25.6M	ImageNet	ImageNet	0.7483	0.7338
Inception_v3	Inception	No	299*299	11.43G	23.8M	ImageNet	ImageNet	0.7401	0.7347
SSD	VGG16	No	300*300	62.77G	26.3M	Voc07+12	Voc07	77.19%	
RefineDet	VGG-16	Support	480*360	123.9G	29.6M	Coco2014	Coco2014	70.14%	
Densebox		Yes	320*320	492M		Private	Private	97.92%	97.50%
Yolo_v3	Yolo	Yes	512*288	53.7G	61.8M	Cityscape	Cityscape	53.7%	53.1%

- DPU deployment
 - **Yes:** the model is successfully deployed on DPU.
 - **Support:** the model is supported but not deployed. Similar model structure is deployed and test successfully.
 - **No:** the model is not supported by DPU right now mainly due to some special operations or layers.

DNNDK – Deep Neural Network Development Kit

> DECENT

- >> Flt32 to Int8 quantization with one line command

> DNNC

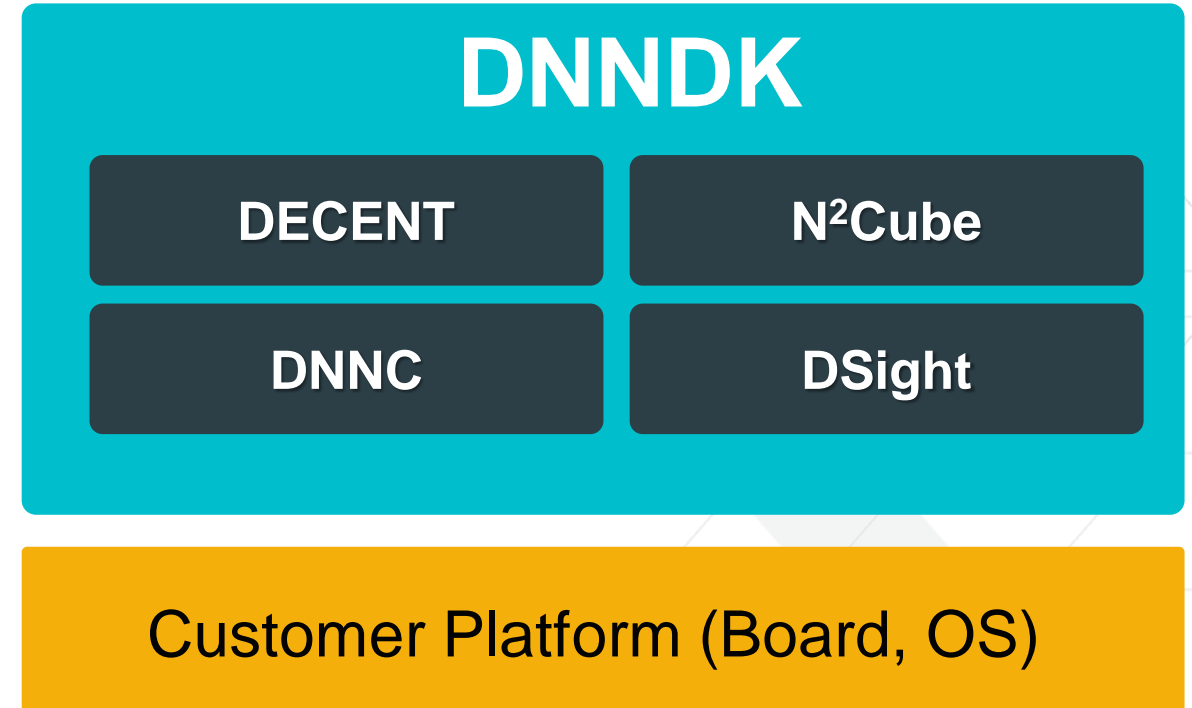
- >> Automatic layer fusion to avoid frequently data read and write

> Runtime N²Cube

- >> Various APIs to facilitate specific application

> Profiler Dsight

- >> Powerful tool as failure analysis and optimization



DNNDK Dev Flow

➤ 01 Model Compression

➤ 02 Model Compilation

➤ 03 Programming

➤ 04 Hybrid Compilation

➤ 05 Execution

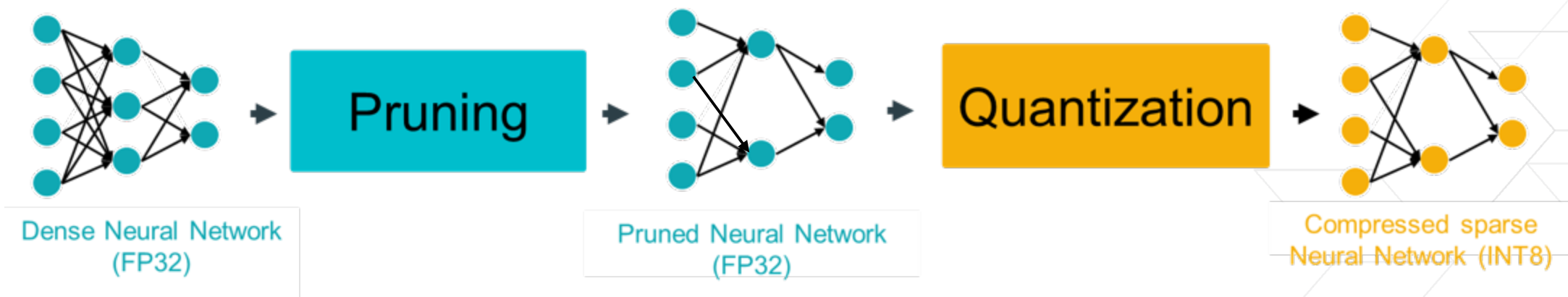
**Five Steps
with DNNDK**

Step 1: Model Compression With DECENT

decent – Deep compression Tool

decent_q – Quantization Tool

decent_p – Pruning Tool

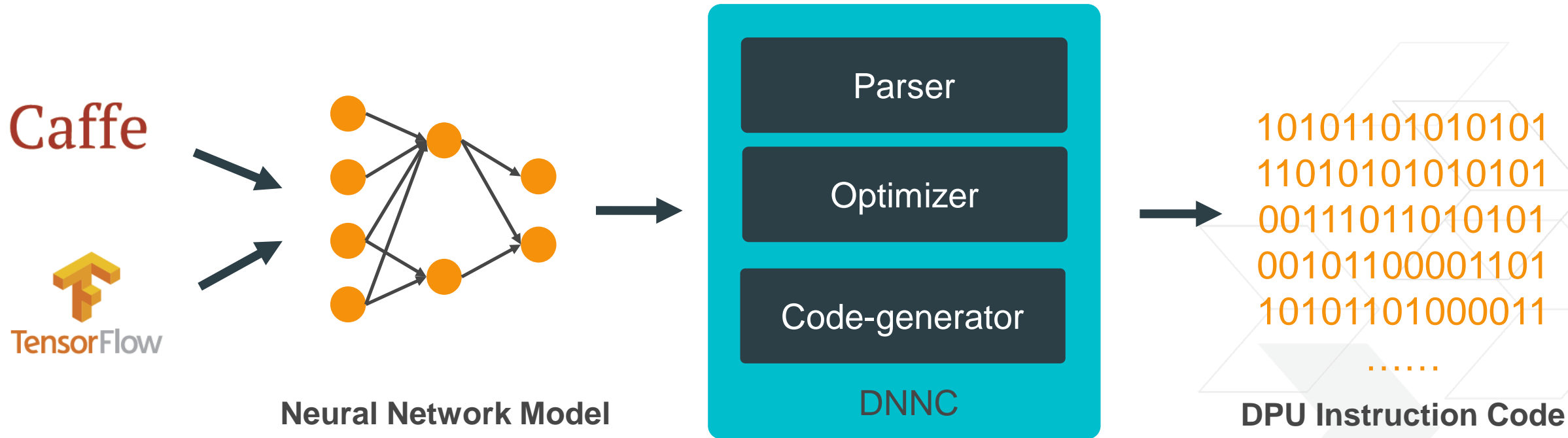


- Consists of two separate tools
 - Quantization Tool
 - Pruning Tool

- Effects
 - Compress model size 5x – 100x
 - Compress running time 1.5x – 10x

- Platform
 - Caffe, Darknet
 - TensorFlow
 - Quantization Tool Beta version
 - Pruning Tool Internal version

Step 2: Model Compilation With DNNC

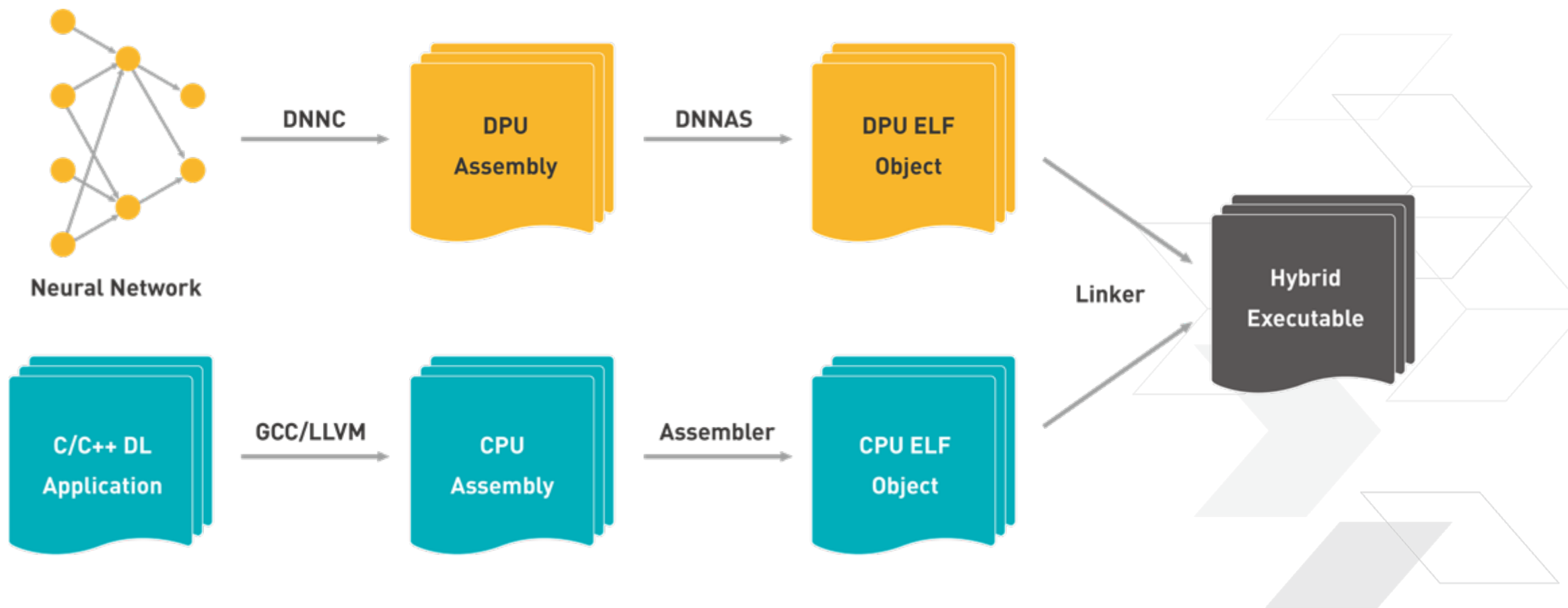


- Programmable tensor-level DPU instruction set
- Compatible for Caffe/TensorFlow frameworks
 - Flexible & scalable for various CNN layers

Step 3: Programming with DNNDK API

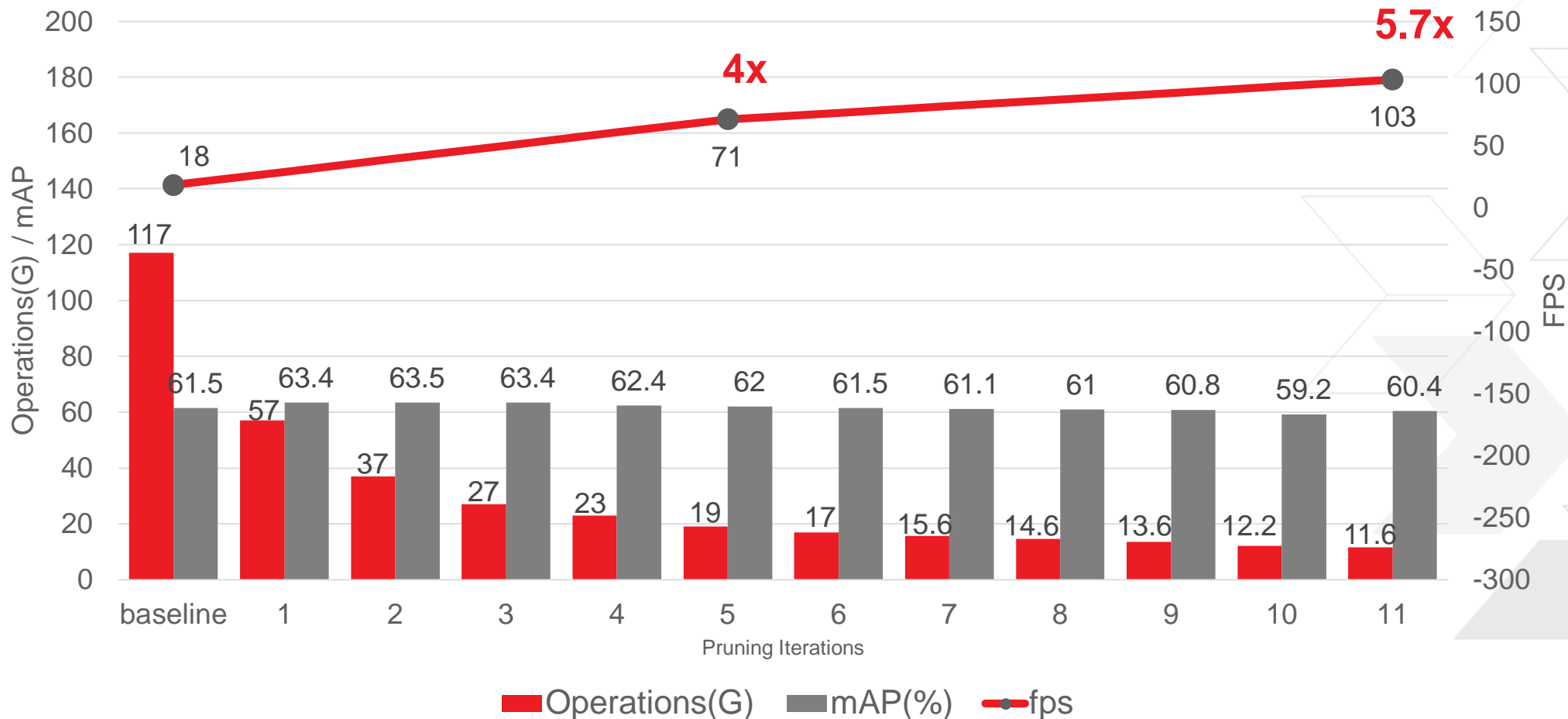
```
1 int main(int argc, char *argv[])
2 {
3     DPUKernel *kernel_conv;
4     DPUKernel *kernel_fc;
5     DPUTask *task_conv;
6     DPUTask *task_fc;
7     char *input_addr;
8     char *output_addr;
9
10    /* DNNDK API to attach to DPU driver */
11    dpuInit();
12
13    /* DNNDK API to create DPU kernels for CONV & FC networks */
14    kernel_conv = dpuLoadKernel("resnet50_conv", 224, 224);
15    kernel_fc = dpuLoadKernel("resnet50_fc", 1, 1);
16
17    /* Create tasks from CONV & FC kernels */
18    task_conv = dpuCreateTask(kernel_conv);
19    task_fc = dpuCreateTask(kernel_fc);
20
21    /* Set input tensor for CONV task and run */
22    input_addr = dpuGetTensorAddress(dpuGetTaskInputTensor(task_conv));
23    setInputImage(Mat &image, input_addr);
24    dpuRunTask(task_conv);
25    output_addr = dpuGetTensorAddress(dpuGetTaskOutputTensor(task_conv));
26
27    /* Run average pooling layer on CPU */
28    run_average_pooling(output_addr);
29
30    /* Set input tensor for FC task and run */
31    input_addr = dpuGetTensorAddress(dpuGetTaskInputTensor(task_fc));
32    setFCInputData(task_fc, input_addr);
33    dpuRunTask(task_fc);
34    output_addr = dpuGetTensorAddress(dpuGetTaskOutputTensor(task_fc));
35
36    /* Display the Classification result from FC task */
37    displayClassificationResult(output_addr);
38
39    /* DNNDK API to destroy DPU tasks/kernels */
40    dpuDestroyTask(task_conv);
41    dpuDestroyTask(task_fc);
42
43    dpuDestroyKernel(kernel_conv);
44    dpuDestroyKernel(kernel_fc);
45
46    /* DNNDK API to detach from DPU driver and free DPU resources */
47    dpuFini();
48
49    return 0;
50 }
```

Step 4: DNNDK Hybrid Compilation Model



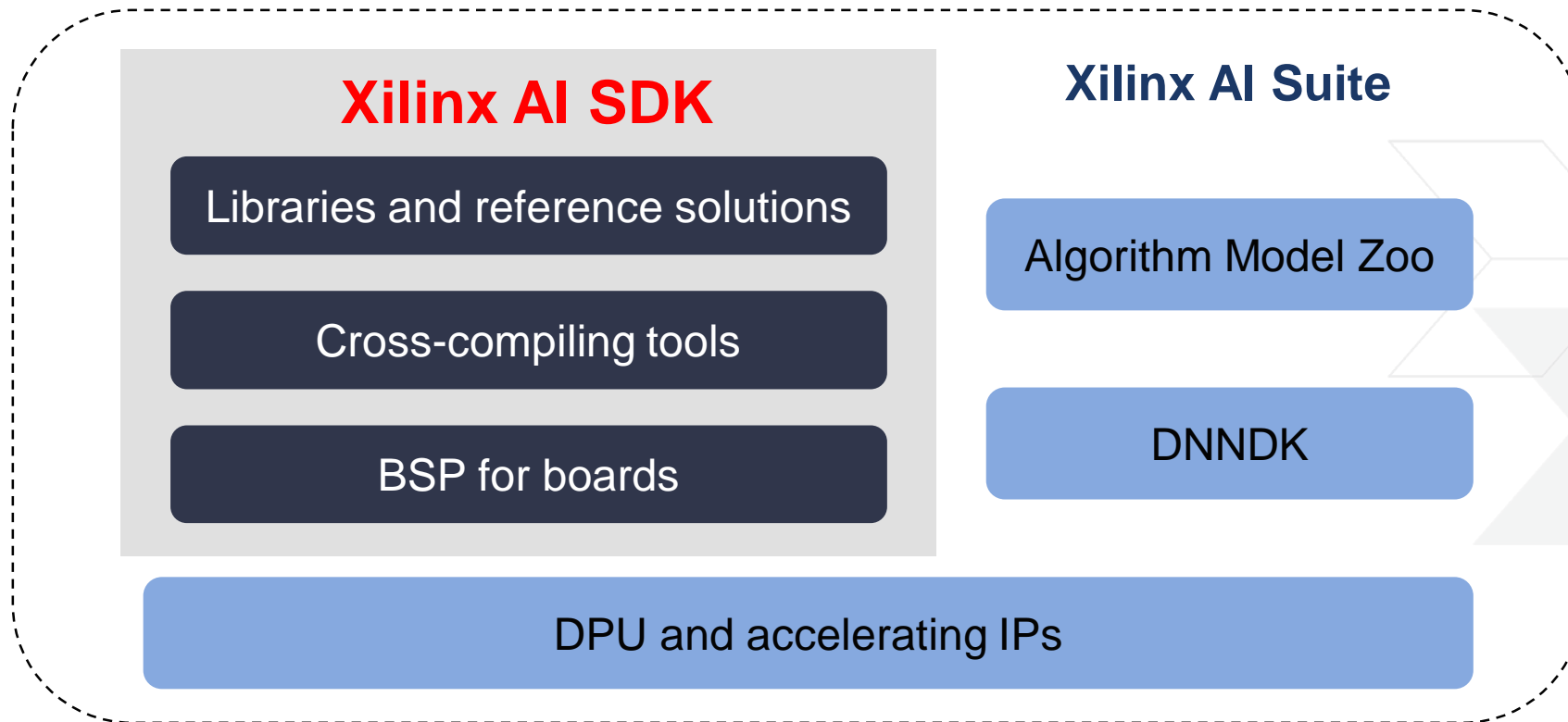
Acceleration on Hardware - SSD

Pruning Speedup on Hardware (2xDPU-4096@ZU9)
VGG_SSD 4 classes detection @Deephi surveillance data

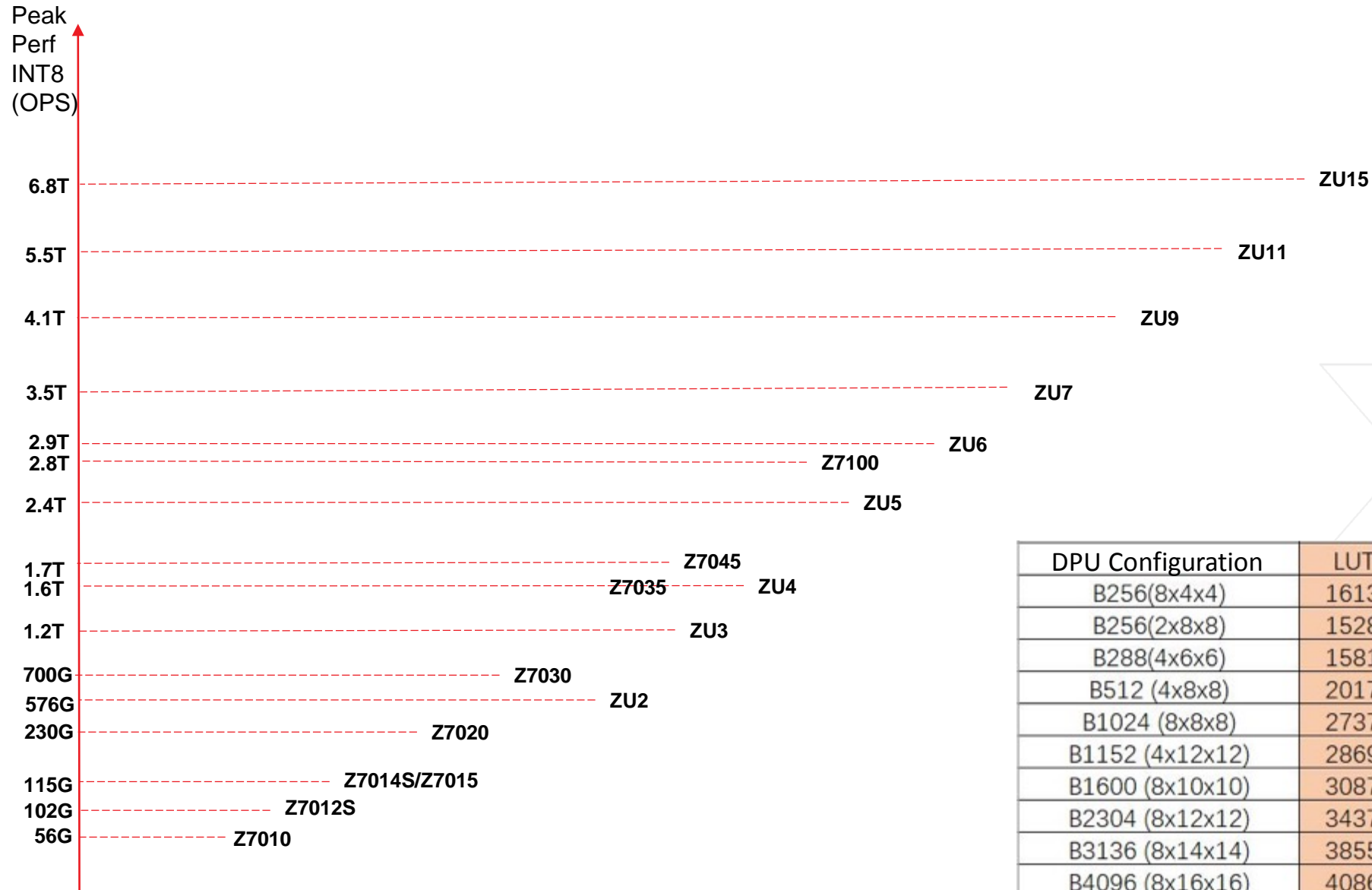


Xilinx AI SDK

- > **Xilinx AI SDK to enable low touch engagements for customers**
 - >> Encourage top customers to use SDK to build applications and solutions
 - >> Only use generic SDK release to support low priority customers



DPU Scalability



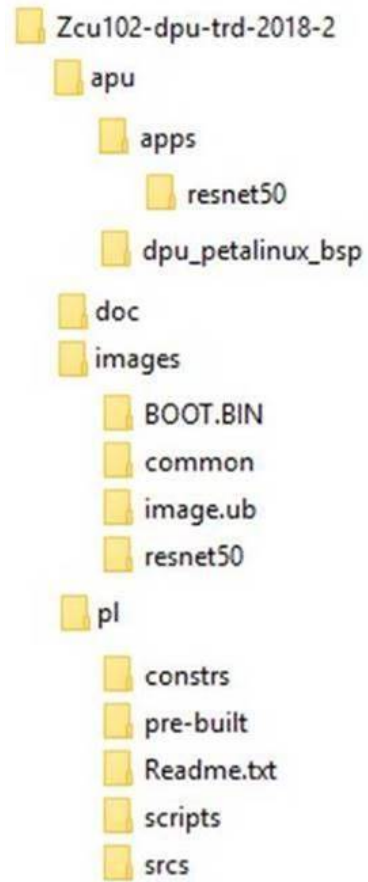
DPU Configuration	LUTs	Registers	BRAM	DSP
B256(8x4x4)	16132	25064	43	66
B256(2x8x8)	15286	22624	53.5	50
B288(4x6x6)	15812	23689	46	62
B512 (4x8x8)	20177	31782	69.5	98
B1024 (8x8x8)	27377	46241	101.5	194
B1152 (4x12x12)	28698	46906	117.5	194
B1600 (8x10x10)	30877	56267	123	282
B2304 (8x12x12)	34379	67481	161.5	386
B3136 (8x14x14)	38555	79867	203.5	506
B4096 (8x16x16)	40865	92630	249.5	642

* B256/288/512/3136 work in progress

DPU IP Integration

The DPU TRD (Targeted Reference Design) has been published in Xilinx.com.

<https://www.xilinx.com/products/design-tools/ai-inference/ai-developer-hub.html#edge>



Customize IP@hw54

dpu_V1301_EU_pack_v1.3.0 (1.3.0)

Documentation IP Location Switch to Defaults

Show disabled ports

Component Name: dpu_V1301_EU_pack_0

Hide	Arch	Summary
	Number of DPUs	1
	Arch of DPU	1152
	Advance	512
		800
		1024
<input checked="" type="checkbox"/>	Independent	1152
	Implementation	1600
	DSP48 Maxim	2304 length 4 [1 - 7]
	DSP48 Use	3136 Maximal
	Ultra-RAM Use per DPU	4096 50

Key Takeaway

- 1 Xilinx offers advanced Edge ML solution
- 2 DNNDK makes neural network deployment easily on FPGA
- 3 DPU IP is ready to download from Xilinx website for free

Adaptable.
Intelligent.