

ABSTRACT

Merchandise Recognition(상품 식별)은 가사/산업용 로봇, 자율주행 자동차 등 다양한 분야에서 필요한 기술이다. 현재 있는 DNN(Deep Neural Network)은 모바일 시스템에서 상품 식별이 어렵다. 기존 DNN system은 높은 accuracy를 갖으나 runtime, area, power가 크다. 따라서 모바일에서 효율적인 DNN accelerator와 algorithm 개발하였다.

Bottleneck layer, Additional layer, ImageDataGenerator 등으로 이루어진 DNN algorithm을 개발하였다. 이를 통해 run-time, data set/size 감소시켰고 accuracy를 높였다.

Energy efficient DNN Accelerator인 zergriss를 개발하였다. ZAC(Zero Aware compression), ZSPE(Zero Skip Processing Element), ZOD(Zero Optimized Dataflow)를 통해 zero를 skip하여 MAC operation을 줄이고 data reuse rate를 높여 DRAM access를 줄였다.

OBJECTIVES

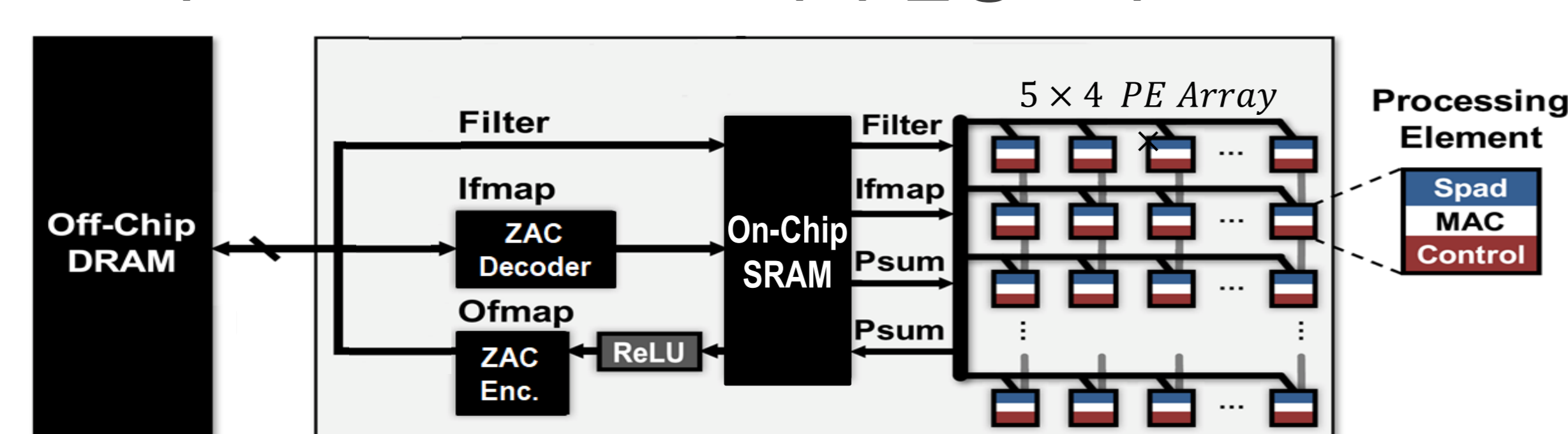
- Software 면에서 상품 구별을 위해 모바일에서 효율적인 DNN algorithm을 개발한다.
 - Data를 limit하고 accuracy를 유지하면서 data set/size, run-time이 감소시킨다.
- Hardware 면에서 상품 구별을 위해 모바일에서 효율적인 DNN accelerator를 개발한다.
 - DRAM access와 MAC operation을 줄인다.

METHODOLOGY

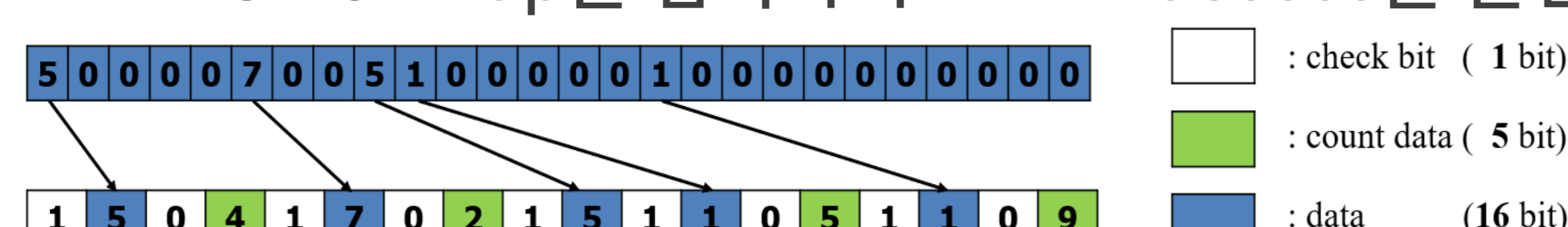
- Bottleneck layer, Additional layer, ImageDataGenerator 등으로 이루어진 DNN algorithm을 개발하였다.
 - python을 이용하여 VGG16, ResNet, MobileNet, InceptionV3 등의 model로 효율을 검증한다.
 - Snack picture 10개 class, 1200 samples로 batch size 16/48, epoch 500, 10GB GPU로 학습/평가한다.



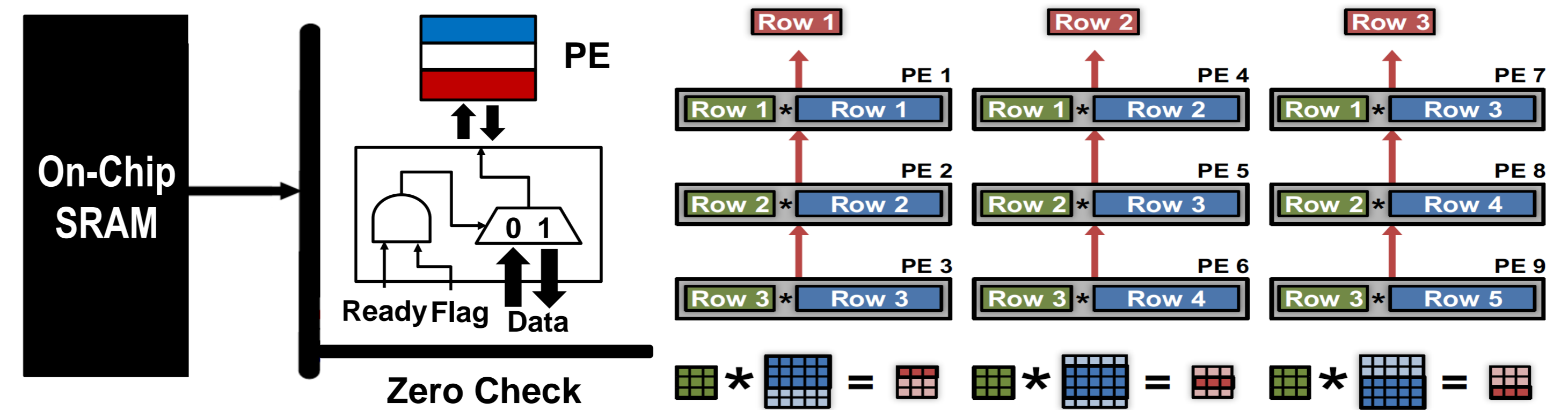
- ZAC, ZSPE, ZOD 등으로 이루어진 DNN accelerator로 Zergriss를 개발하였다.
 - 구조는 DRAM, SRAM, PE, Encoder, Decoder이다.
 - 1차 : python으로 DNN accelerator를 Lenet-5에서 검증한다.
 - 2차 : verilog으로 작성 후 modelsim으로 검증한다.
 - 3차 : vivado로 FPGA에서 검증한다.



- ZAC : OFMap을 압축하여 DRAM access를 줄인다.

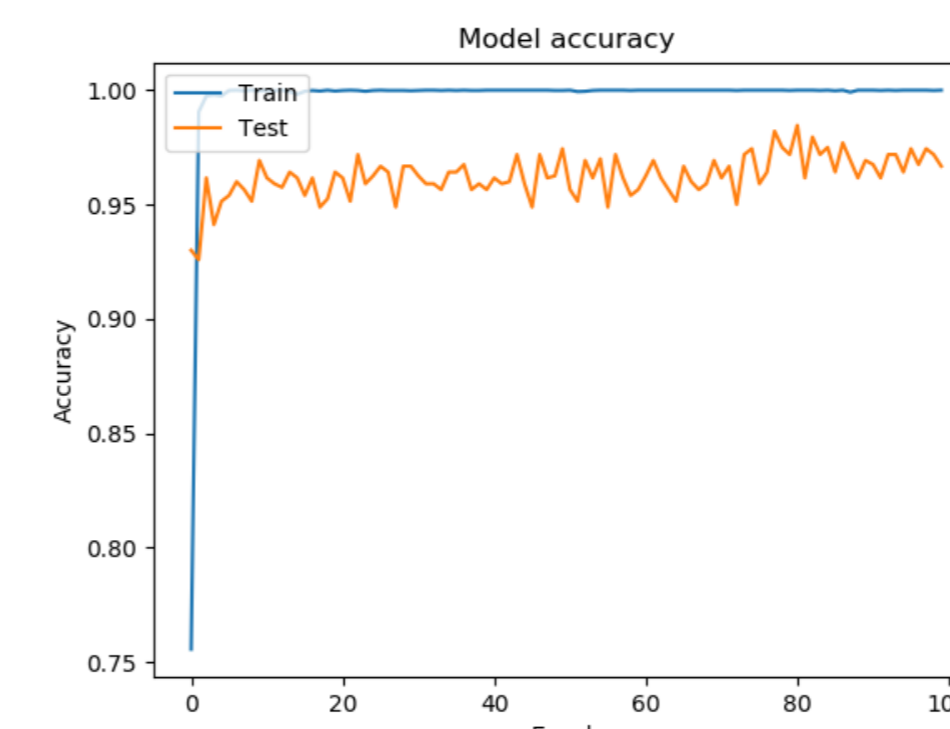


- ZSPE : IFMap(Image Feature Map)와 Filter data 중 Zero를 aware/skip하여 MAC operation을 줄인다.
- ZOD : RS dataflow로 IFMap과 filter를 reuse하여 DRAM access를 줄인다.

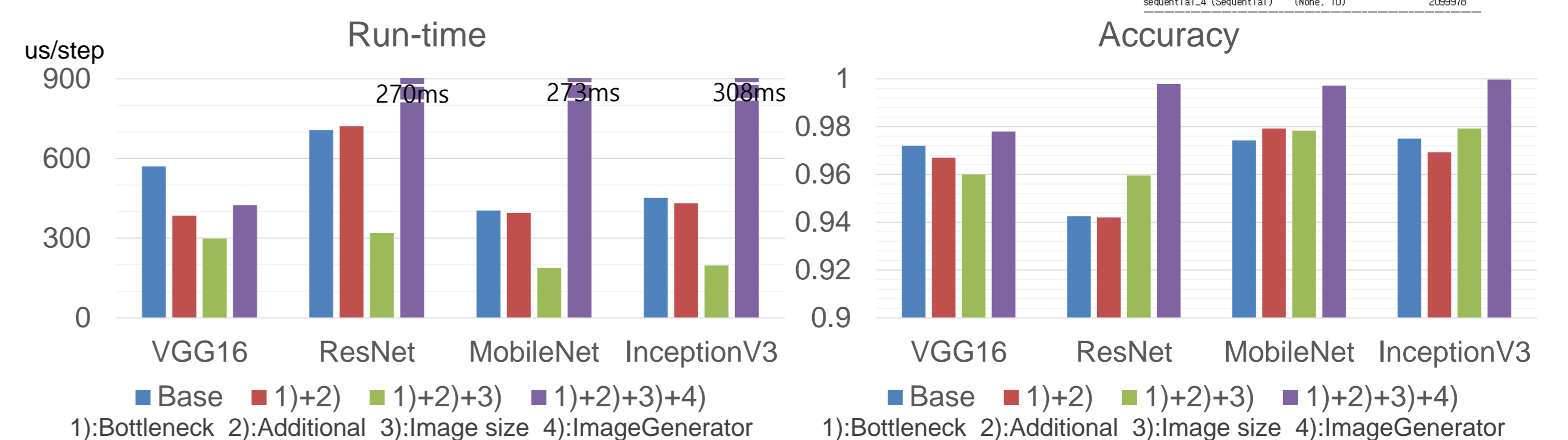


RESULTS

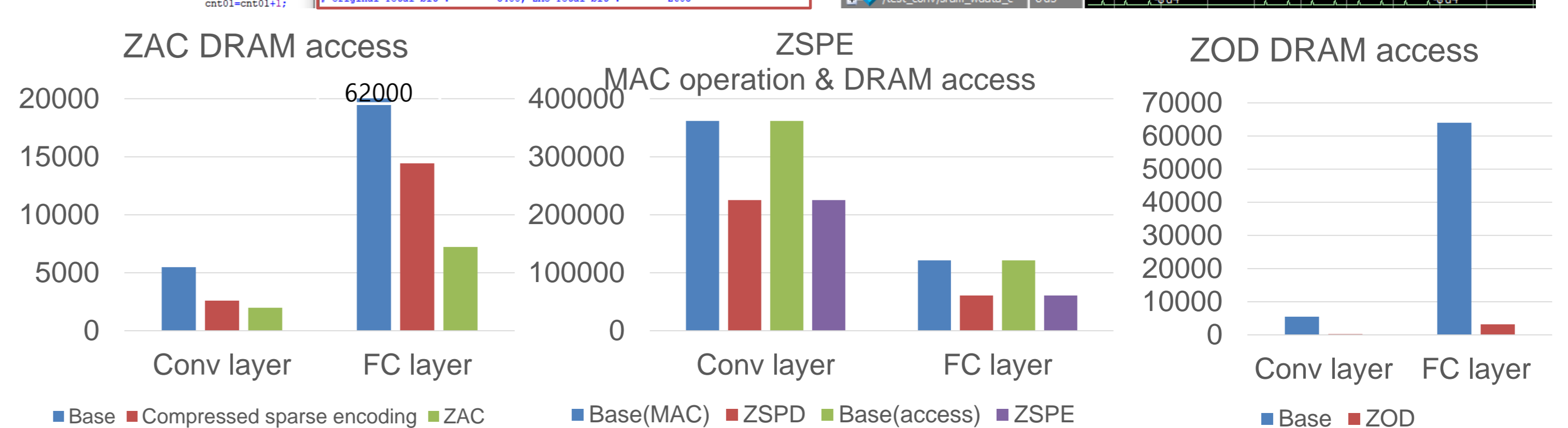
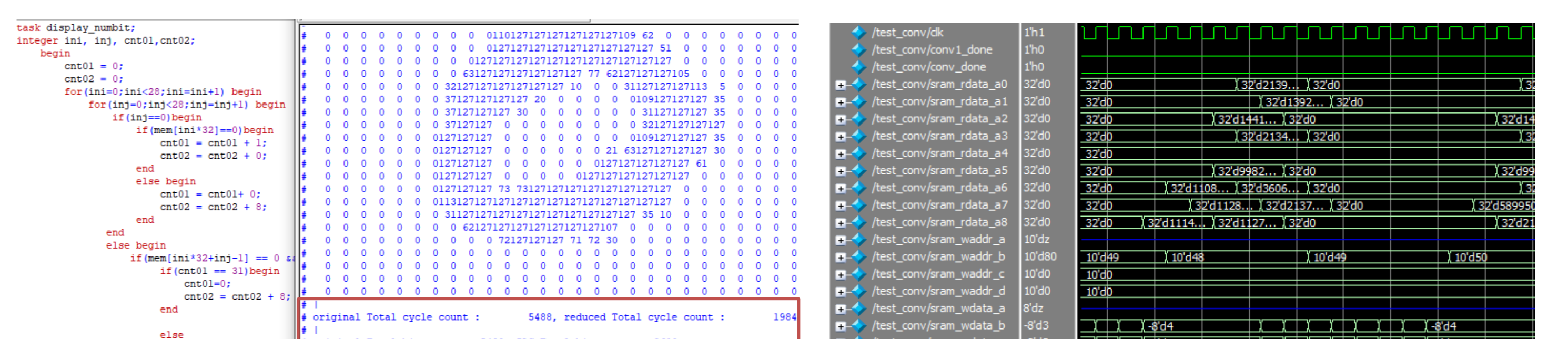
- 한 상품의 training data를 100개로 limit하고 run time을 53%로 줄이고 accuracy를 2.8% 향상시켰다.



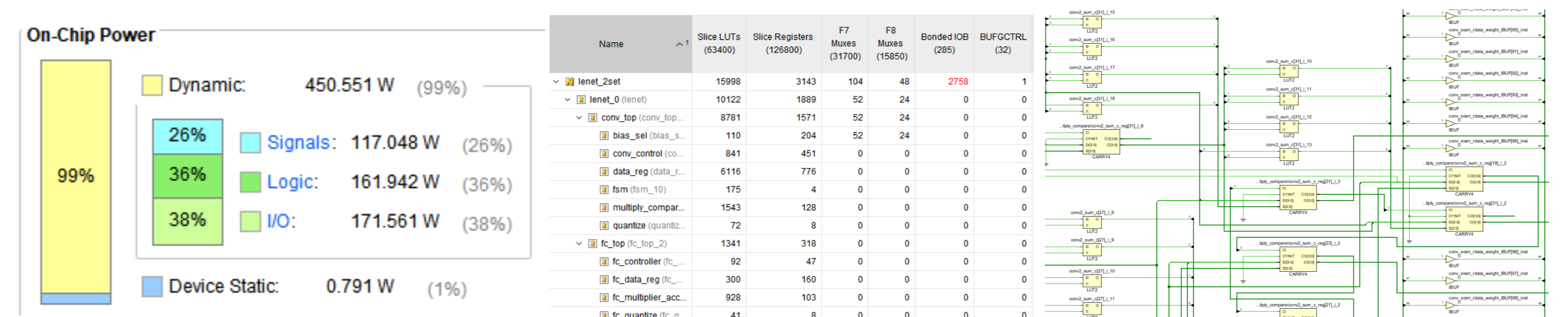
Layer	Group	Output Shape	Para #
lenet_1	lenet_1	(None, 10, 10, 3)	0
lenet_2	lenet_2	(None, 10, 10, 16)	176
lenet_3	lenet_3	(None, 10, 10, 16)	256
lenet_4	lenet_4	(None, 10, 10, 16)	256
lenet_5	lenet_5	(None, 10, 10, 16)	256
lenet_6	lenet_6	(None, 10, 10, 16)	256
lenet_7	lenet_7	(None, 10, 10, 16)	256
lenet_8	lenet_8	(None, 10, 10, 16)	256
lenet_9	lenet_9	(None, 10, 10, 16)	256
lenet_10	lenet_10	(None, 10, 10, 16)	256
lenet_11	lenet_11	(None, 10, 10, 16)	256
lenet_12	lenet_12	(None, 10, 10, 16)	256
lenet_13	lenet_13	(None, 10, 10, 16)	256
lenet_14	lenet_14	(None, 10, 10, 16)	256
lenet_15	lenet_15	(None, 10, 10, 16)	256
lenet_16	lenet_16	(None, 10, 10, 16)	256
lenet_17	lenet_17	(None, 10, 10, 16)	256
lenet_18	lenet_18	(None, 10, 10, 16)	256
lenet_19	lenet_19	(None, 10, 10, 16)	256
lenet_20	lenet_20	(None, 10, 10, 16)	256
lenet_21	lenet_21	(None, 10, 10, 16)	256
lenet_22	lenet_22	(None, 10, 10, 16)	256
lenet_23	lenet_23	(None, 10, 10, 16)	256
lenet_24	lenet_24	(None, 10, 10, 16)	256
lenet_25	lenet_25	(None, 10, 10, 16)	256
lenet_26	lenet_26	(None, 10, 10, 16)	256
lenet_27	lenet_27	(None, 10, 10, 16)	256
lenet_28	lenet_28	(None, 10, 10, 16)	256
lenet_29	lenet_29	(None, 10, 10, 16)	256
lenet_30	lenet_30	(None, 10, 10, 16)	256
lenet_31	lenet_31	(None, 10, 10, 16)	256
lenet_32	lenet_32	(None, 10, 10, 16)	256
lenet_33	lenet_33	(None, 10, 10, 16)	256
lenet_34	lenet_34	(None, 10, 10, 16)	256
lenet_35	lenet_35	(None, 10, 10, 16)	256
lenet_36	lenet_36	(None, 10, 10, 16)	256
lenet_37	lenet_37	(None, 10, 10, 16)	256
lenet_38	lenet_38	(None, 10, 10, 16)	256
lenet_39	lenet_39	(None, 10, 10, 16)	256
lenet_40	lenet_40	(None, 10, 10, 16)	256
lenet_41	lenet_41	(None, 10, 10, 16)	256
lenet_42	lenet_42	(None, 10, 10, 16)	256
lenet_43	lenet_43	(None, 10, 10, 16)	256
lenet_44	lenet_44	(None, 10, 10, 16)	256
lenet_45	lenet_45	(None, 10, 10, 16)	256
lenet_46	lenet_46	(None, 10, 10, 16)	256
lenet_47	lenet_47	(None, 10, 10, 16)	256
lenet_48	lenet_48	(None, 10, 10, 16)	256
lenet_49	lenet_49	(None, 10, 10, 16)	256
lenet_50	lenet_50	(None, 10, 10, 16)	256
lenet_51	lenet_51	(None, 10, 10, 16)	256
lenet_52	lenet_52	(None, 10, 10, 16)	256
lenet_53	lenet_53	(None, 10, 10, 16)	256
lenet_54	lenet_54	(None, 10, 10, 16)	256
lenet_55	lenet_55	(None, 10, 10, 16)	256
lenet_56	lenet_56	(None, 10, 10, 16)	256
lenet_57	lenet_57	(None, 10, 10, 16)	256
lenet_58	lenet_58	(None, 10, 10, 16)	256
lenet_59	lenet_59	(None, 10, 10, 16)	256
lenet_60	lenet_60	(None, 10, 10, 16)	256
lenet_61	lenet_61	(None, 10, 10, 16)	256
lenet_62	lenet_62	(None, 10, 10, 16)	256
lenet_63	lenet_63	(None, 10, 10, 16)	256
lenet_64	lenet_64	(None, 10, 10, 16)	256
lenet_65	lenet_65	(None, 10, 10, 16)	256
lenet_66	lenet_66	(None, 10, 10, 16)	256
lenet_67	lenet_67	(None, 10, 10, 16)	256
lenet_68	lenet_68	(None, 10, 10, 16)	256
lenet_69	lenet_69	(None, 10, 10, 16)	256
lenet_70	lenet_70	(None, 10, 10, 16)	256
lenet_71	lenet_71	(None, 10, 10, 16)	256
lenet_72	lenet_72	(None, 10, 10, 16)	256
lenet_73	lenet_73	(None, 10, 10, 16)	256
lenet_74	lenet_74	(None, 10, 10, 16)	256
lenet_75	lenet_75	(None, 10, 10, 16)	256
lenet_76	lenet_76	(None, 10, 10, 16)	256
lenet_77	lenet_77	(None, 10, 10, 16)	256
lenet_78	lenet_78	(None, 10, 10, 16)	256
lenet_79	lenet_79	(None, 10, 10, 16)	256
lenet_80	lenet_80	(None, 10, 10, 16)	256
lenet_81	lenet_81	(None, 10, 10, 16)	256
lenet_82	lenet_82	(None, 10, 10, 16)	256
lenet_83	lenet_83	(None, 10, 10, 16)	256
lenet_84	lenet_84	(None, 10, 10, 16)	256
lenet_85	lenet_85	(None, 10, 10, 16)	256
lenet_86	lenet_86	(None, 10, 10, 16)	256
lenet_87	lenet_87	(None, 10, 10, 16)	256
lenet_88	lenet_88	(None, 10, 10, 16)	256
lenet_89	lenet_89	(None, 10, 10, 16)	256
lenet_90	lenet_90	(None, 10, 10, 16)	256
lenet_91	lenet_91	(None, 10, 10, 16)	256
lenet_92	lenet_92	(None, 10, 10, 16)	256
lenet_93	lenet_93	(None, 10, 10, 16)	256
lenet_94	lenet_94	(None, 10, 10, 16)	256
lenet_95	lenet_95	(None, 10, 10, 16)	256
lenet_96	lenet_96	(None, 10, 10, 16)	256
lenet_97	lenet_97	(None, 10, 10, 16)	256
lenet_98	lenet_98	(None, 10, 10, 16)	256
lenet_99	lenet_99	(None, 10, 10, 16)	256
lenet_100	lenet_100	(None, 10, 10, 16)	256



- Modelsim으로 Lenet5에서 bit, cycle을 비교하였다.
 - ZAC로 DRAM access를 기존보다 76% 줄였다.
 - ZSPE : MAC operation을 40.8%, DRAM access를 40.7% 줄였다.
 - ZOD : DRAM access를 Lenet5 대비 95% 줄 것이다. (testbench에서 ZOD 사용 시 reuse 횟수를 통해 예측한 값 / 미 구현)



- Vivado로 board i/o pins의 한계로 synthesis까지 진행하여 Lenet5에서 zergriss 효율을 검증하였다.



CONCLUSIONS

개발된 DNN algorithm은 기존보다 적은 data에서 runtime은 53% 줄였고 accuracy는 2.8% 향상시켰다. 개발된 DNN accelerator은 기존보다 DRAM access를 85.8%, MAC operation을 40.8% 줄였다. 추후 algorithm으로는 pruning, skipping 관련 기술과 ZOD 구현하여 다른 모델을 FPGA에서 추가 검증한다.