Week 2: A Big Picture?

### Last Week:

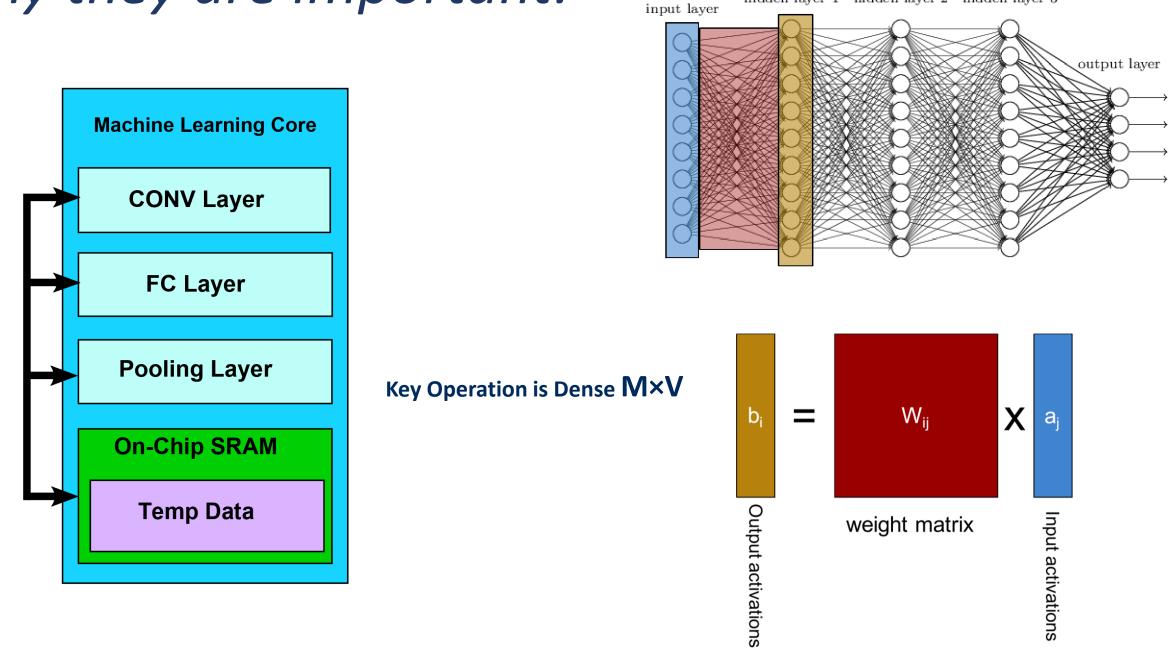
#### **CONV** Layer

**FC Layer** 

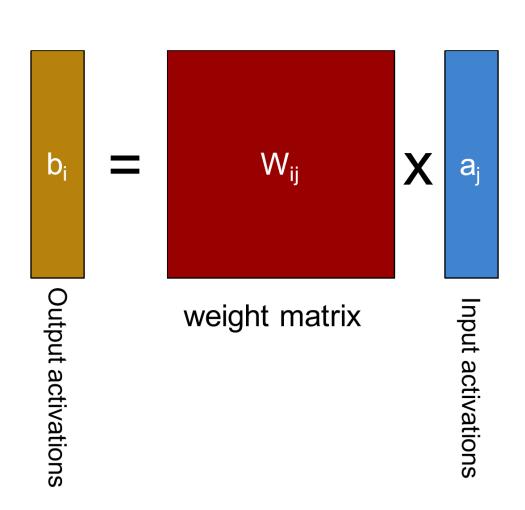
#### **Questions:**

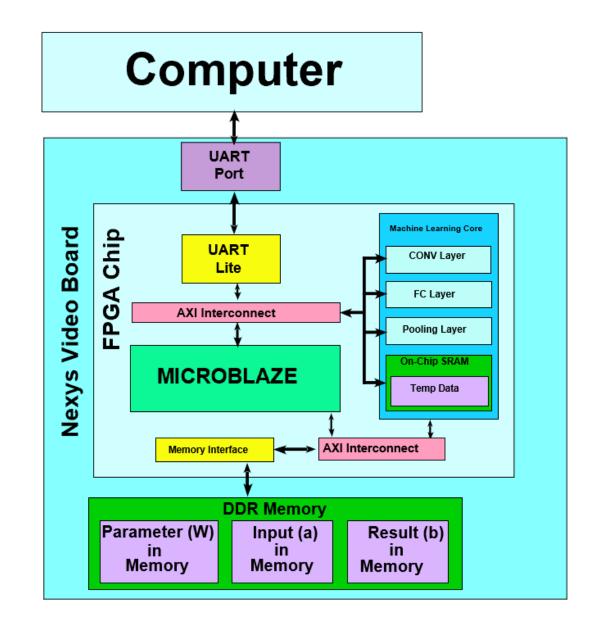
- 1. Why are they important?
- 2. What are we doing?
- 3. How are we going to use them?
- 4. Give me 3 other solutions? (Can we do better?)

# Why they are important:

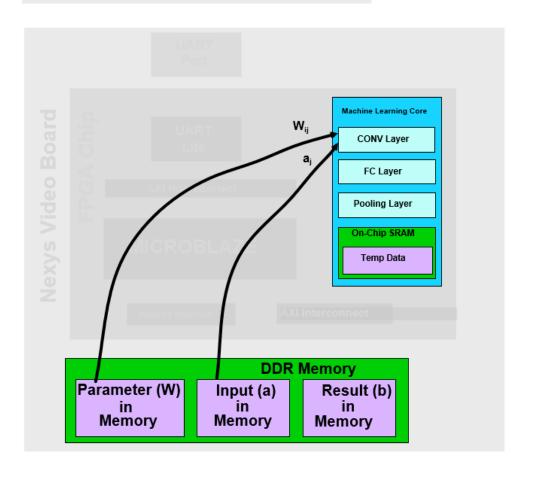


hidden layer 1 hidden layer 2 hidden layer 3

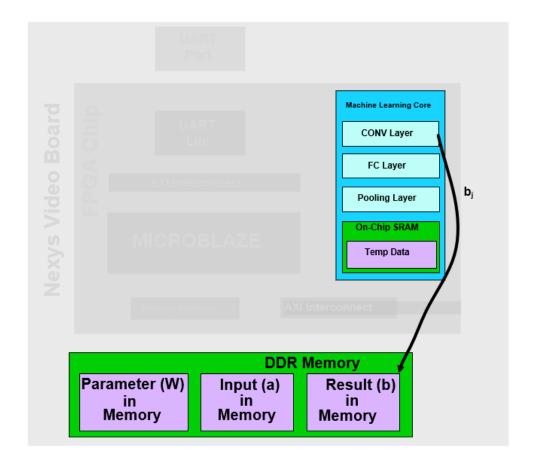




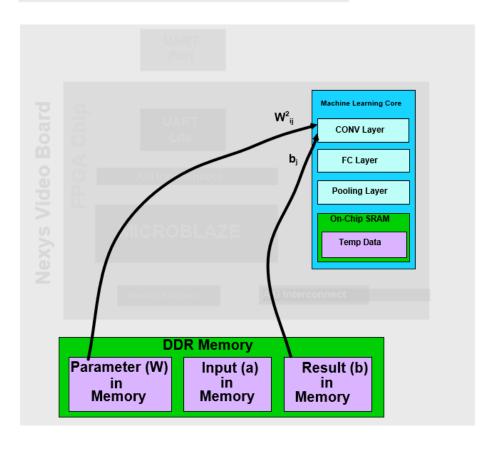
### Computer



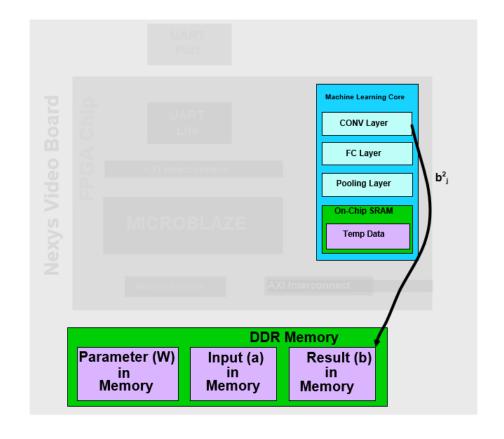
### Computer



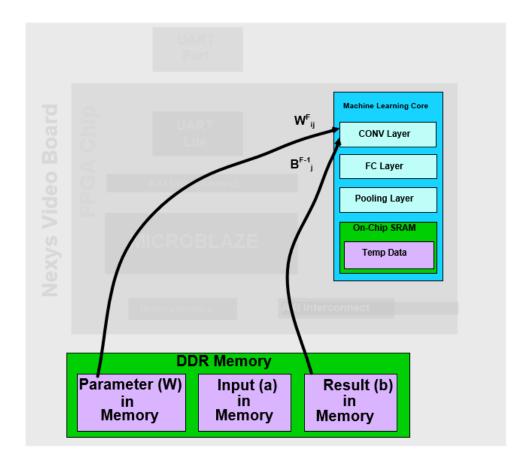
Computer



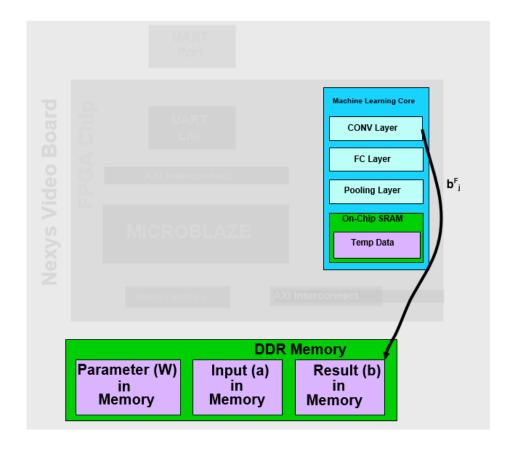
#### Computer



Computer



### Computer



# How we are using them:



Synthesis 'C/C++' to RTL.



\/\\\A\\) Create blocks, Gate level design



Program Microprocessor such that it can interface with different IPs.

# How we are using them:

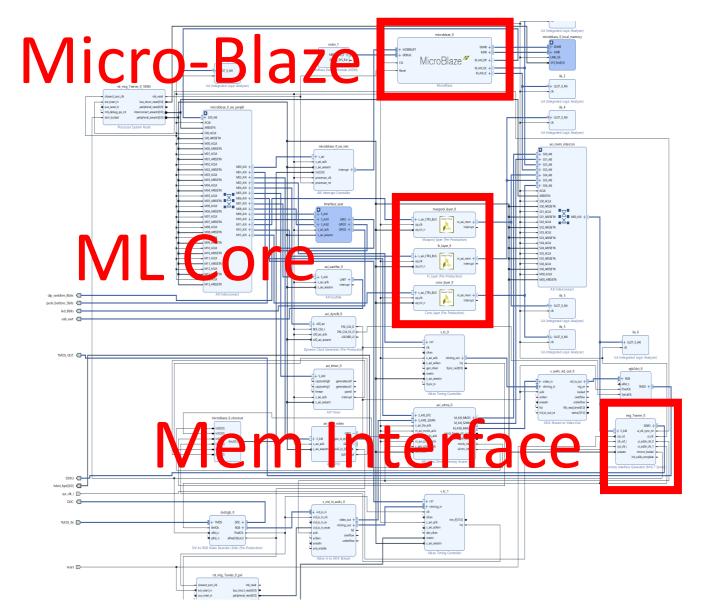




Diagram × Address Editor ×					
Q 🗶 🖨					
Cell	Slave Interface	Slave Segment	Offset Address	Range	High Address
✓ ▼ maxpool_layer_0					
∨ ■ Data_m_axi_mem (32 address bits : 4G)					
mig_7series_0	S_AXI	memaddr	0x8000_0000	512M *	0x9FFF_FFFF
✓ ≢ conv_layer_0					
✓ ■ Data_m_axi_mem (32 address bits : 4G)					
mig_7series_0	S_AXI	memaddr	0x8000_0000	512M -	0x9FFF_FFFF
✓ 🕴 fc_layer_0					
✓ ■ Data_m_axi_mem (32 address bits : 4G)					
mig_7series_0	S_AXI	memaddr	0x8000_0000	512M -	0x9FFF_FFFF
∨ ₹ axi_vdma_0					
✓ ■ Data_MM2S (32 address bits : 4G)					
mig_7series_0	S_AXI	memaddr	0x8000_0000	512M -	0x9FFF_FFFF
✓ ■ Data_S2MM (32 address bits : 4G)					
mig_7series_0	S_AXI	memaddr	0x8000_0000	512M -	0x9FFF_FFFF
∨ ₹ microblaze_0					
✓ ■ Data (32 address bits : 4G)					
axi_dynclk_0	s00_axi	reg0	0x44A2_0000	64K *	0x44A2_FFFF
Interface_user/axi_gpio_SW	S_AXI	Reg	0x4002_0000	64K *	0x4002_FFFF
Interface_user/axi_gpio_Push	S_AXI	Reg	0x4003_0000	64K *	0x4003_FFFF
Interface_user/axi_gpio_LED	S_AXI	Reg	0x4001_0000	64K *	0x4001_FFFF
axi_gpio_video	S_AXI	Reg	0x4000_0000	64K *	0x4000_FFFF
axi_timer_0	S_AXI	Reg	0x41C0_0000	64K -	0x41C0_FFFF
axi_uartlite_0	S_AXI	Reg	0x4060_0000	64K -	0x4060_FFFF
	CAMILITE		0-4440 0000	CAV	04440 PEPP
conv_layer_0	s_axi_CTRL_BUS	Reg	0x44A4_0000	64K -	0x44A4_FFFF
,					
fc_layer_0	s_axi_CTRL_BUS	Reg	0x44A5_0000	64K *	0x44A5_FFFF
maxpool_layer_0	s_axi_CTRL_BUS	Reg	0x44A6_0000	64K *	0x44A6_FFFF
microblaze_0_axi_intc	s_axı	Reg	0X4120_0000	64K *	0X4120_FFFF
mig_7series_0	S_AXI	memaddr	0x8000_0000	512M -	0x9FFF_FFFF
v_tc_0	ctrl	Reg	0x44A1_0000	64K -	0x44A1_FFFF
□ v_tc_1	ctrl	Reg	0x44A3_0000	64K *	0x44A3_FFFF
✓ ■ Instruction (32 address bits : 4G)					
microblaze_0_local_memory/ilmb_bram_if_cntlr	SLMB	Mem	0x0000_0000	32K -	0x0000_7FFF
mig_7series_0	S_AXI	memaddr	0x8000_0000	512M *	0x9FFF_FFFF

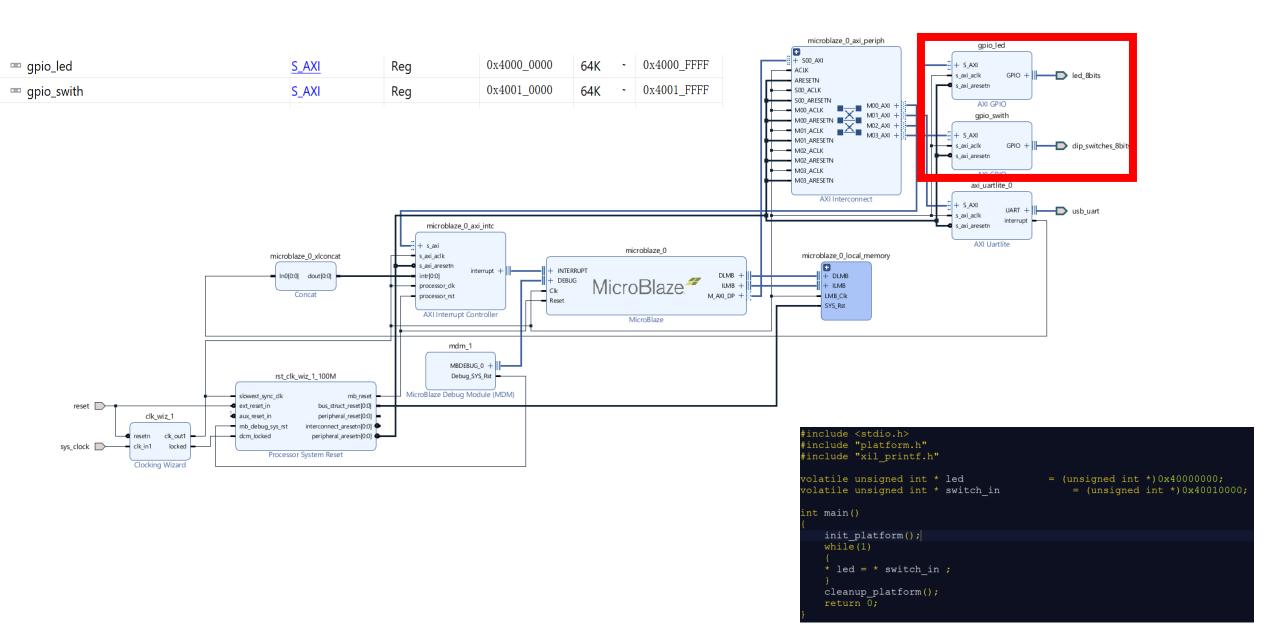
# How we are using them:



```
void hw hls core processing ()
   printf("start hls core processing \n");
   //layer 1 - conv 1
      hw conv layer(2415919104, 2432696320, 1, 28, 28, 20, 24, 24, 1, 5);
      input offset 0x90000000, output offset 0x91000000
   //layer 2 -pool 1
   hw maxpool layer(2432696320, 2449573736, 20, 24, 24, 20, 12, 12, 2, 2);
   input offset 0x91000000, output offset 0x92018768 (conv2 bias + weight = 50
   + 25000 = 25050 = 0x18768
   //layer 3 - conv 2
      hw conv layer(2449473536, 2466250752, 20, 12, 12, 50, 8, 8, 1, 5);
      input offset 0x92000000 , output offset 0x93000000
   //layer 4 -pool 2
   hw maxpool layer(2466250752, 2484629968, 50, 8, 8, 50, 4, 4, 2, 2);
   input offset 0x93000000, output offset 0x941871D0 (fcl bias + weight =500 +
   400000 = 400500 = 0x1871D0
   //layer 5 -fc 1
        hw fc layer(2483027968, 2499825224, 800, 500, 1);
        input offset 0x94000000 , output offset 0x95004E48 (fc2 bias + weight
        =10 + 5000 = 5010 = 0x4E48
   //laver 6 -fc 2
        hw fc layer(2499805184, 2516582400, 500, 10, 0);
        input offset 0x95000000, output offset 0x96000000
   printf("done hls core processing \n");
```

```
include <algorithm>
                                                                     Vivado™ HLS
include "fc layer.h"
oid fc layer(float * mem
             int input offset,
            int output offset,
            const int batch size,
             const int num inputs,
             const int num outputs,
             const int enable relu)
pragma HLS INTERFACE m axi port=mem depth=2147483648
pragma HLS INTERFACE s axilite port=input offset bundle=CTRL BUS
pragma HLS INTERFACE s_axilite port=output_offset bundle=CTRL_BUS
pragma HLS INTERFACE s axilite port=batch size bundle=CTRL BUS
pragma HLS INTERFACE s axilite port=num inputs bundle=CTRL BUS
pragma HLS INTERFACE s axilite port=num outputs bundle=CTRL BUS
pragma HLS INTERFACE s_axilite port=enable_relu bundle=CTRL_BUS
pragma HLS INTERFACE s axilite port=return bundle=CTRL BUS
 const int num_weights = num_inputs*num_outputs;
 const int num biases = num outputs;
 for (int b = 0; b < batch size; b++) {</pre>
   for (int o = 0; o < num outputs; o++) {
     // Set bias
     float output element = mem[input offset/sizeof(float) + num weights + o];
     for (int i = 0; i < num inputs; i = i+1) {
 #pragma HLS PIPELINE
       float input element = mem[input offset/sizeof(float) + num weights +
      num biases + b*num inputs+i];
           float weight element = mem[input offset/sizeof(float) + o*num inputs+i];
      // if (weight element != 0.0f) {
               output element += input element * weight element;
      //float input element = mem[input offset/sizeof(float) + num weights +
      num biases + b*num inputs+i];
      //float weight element = mem[input offset/sizeof(float) + o*num inputs+i];
     // Compute activation
     if (enable relu)
       mem[output offset/sizeof(float) + b*num outputs+o] = std::max(0.0f,
       output element);
        mem[output offset/sizeof(float) + b*num outputs+o] = output element;
```

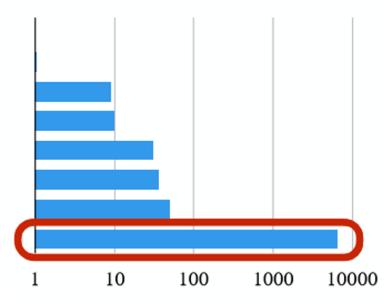
### Next Week Tutorial:

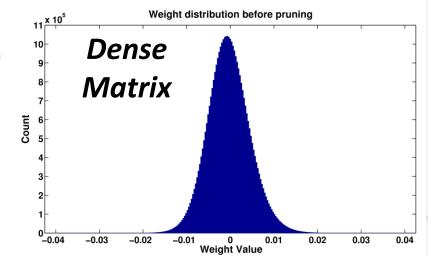


### Can we do better? - Deep Compression.

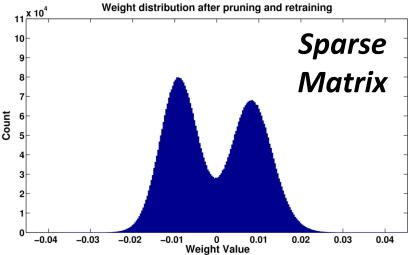
Relative Energy Cost

Operation	Energy [pJ]	Relative Cost
32 bit int ADD	0.1	1
32 bit float ADD	0.9	9
32 bit Register File	1	10
32 bit int MULT	3.1	31
32 bit float MULT	3.7	37
32 bit SRAM Cache	5	50
32 bit DRAM Memory	640	6400



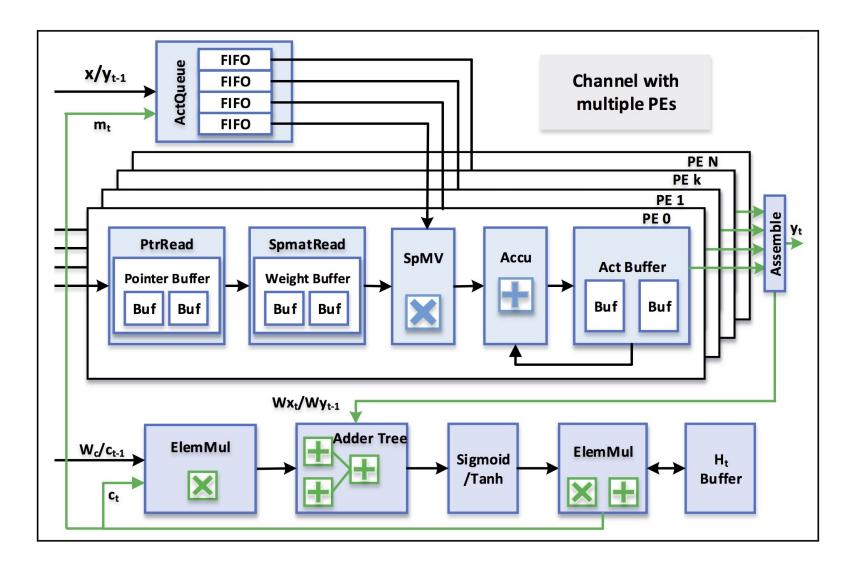






FROM: S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. Horowitz, W. Dally, "EIE: Efficient Inference Engine on Compressed Deep Neural Network", ISCA'16.

### Can we do better? -Better Architecture.



From: S. Han, J. Kang, H. Mao, Y. Li, D. Xie, H. Luo, Y. Wang, H. Yang, W.J. Dally "ESE: Efficient Speech Recognition Engine for Compressed LSTM", to appear at FPGA'17

- 1. Get something working.
- 2. Novel Architecture/Algorithm.

Thank You!

Have a nice weekend!