

9

Hardware Architecture for Deep Learning - CS6490, Spring 2024-25.
 Dept. of CSE, IIT Hyderabad
Quiz-1, Set-B

Time: 10 minutes

Total marks: 10

Name: Pandriang Aditya Sagaram

Roll number: EE22 B TECH 11039

1. The depth (z-dimension) of a kernel in a CONV layer is same as (choose correct ones): [1]

- a. The number of input channels
- b. The number of output channels
- c. The number of columns of the output
- d. The number of rows of the input



2. The number of kernels/filters in a CONV layer of a CNN is same as (choose correct ones): [1]

- a. The number of input channels
- b. The number of output channels
- c. The number of columns of the output
- d. The number of rows of the input

3. Consider a CONV layer of a CNN with following characteristic: Input: $192 \times 192 \times 7$, #filters=48, filter size = 3x3, stride=1. Find the dimensions of the output. [2]

$$\text{Stride} = 1 \Rightarrow \text{Output size} = 192 - 3 + 1 = 190 \\ (\text{length/width})$$

~~Depth = 7x48
Channels~~

$$\text{Output channels} = 7 \\ \text{per filter}$$

$$\begin{array}{r} 48 \\ \times 2 \\ \hline 96 \\ \hline 336 \end{array}$$

∴ with 48 filters, dimensions are $190 \times 190 \times 336$

4. Consider an FC layer of a CNN with following characteristics: Input: $7 \times 7 \times 96$, #filters=2048. Find the dimensions of the filter and the output. [2]

Filter : 7×7 each (since it is FC) and 96 channels

$$\text{Output Dim : } 7 \times 7 - 7 + 1 = 1 \times 1$$

$$\rightarrow 1 \times 1 \times 2048$$

2

5. Perform 2-D convolution of the given 4×4 input with the 2×2 kernel. Assume stride=2. [4]

Input:

5	7	2	1
6	1	2	5
1	2	0	2
11	3	2	9

$$\text{Output size} = \left(\frac{4}{2} + \frac{2}{2} - 1 \right) \\ (\text{length or width}) \\ = 2$$

Kernel:

-1	2
-3	4

$\therefore 2 \times 2 \text{ matrix}$

$$\text{out}[0][0] = 5(-1) + 7(2) + 6(-3) + 1(4) \\ = -5 + 14 - 18 + 4 \\ = -5$$

$$\text{out}[0][1] = 2(-1) + 1(2) + 2(-3) + 5(4) \\ = -2 + 2 - 6 + 20 \\ = 14$$

$$\text{out}[1][0] = 1(-1) + 2(2) + 11(-3) + 3(4) \\ = -1 + 4 - 33 + 12 \\ = -18$$

$$\text{out}[1][1] = 0 + 2(2) + (-3)(2) + (4)(4) \\ = 4 - 6 + 36 \\ = 34$$

$$\therefore \text{output} = \begin{bmatrix} -5 & 14 \\ -18 & 34 \end{bmatrix}$$

4



(a)

Hardware Architecture for Deep Learning - CS6490. Spring 2024-25.
Dept. of CSE, IIT Hyderabad

Quiz-2, Set-A

Time: 10 minutes

Total marks: 10

Name: Pandiangi Aditya Siram

Roll number: EE22BTECH11039

Choose all correct choices for the MCQ.

1. A systolic array enables: [2]

- a. Parallel processing
- b. Data reuse
- c. Run-length compression
- d. Simpler design

2. Which dataflow is depicted by the given code snippet for a 1-D convolution: [1]

```
for (i=0; i<8; i++){  
    out[i] = 0;  
    for (j=0; j<3; j++)  
        out[i] = out[i] + in[i+j]*w[j];  
}
```

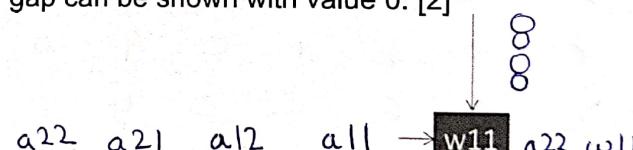
out, in, wt represent the output, input, and weight matrices, respectively.

- a. Output stationary
- b. Input stationary
- c. Weight stationary
- d. Row stationary

3. Given the systolic design as shown, mark the values to be fed to horizontal and vertical inputs to perform the convolution operation. Any cycle gap can be shown with value 0. [2]

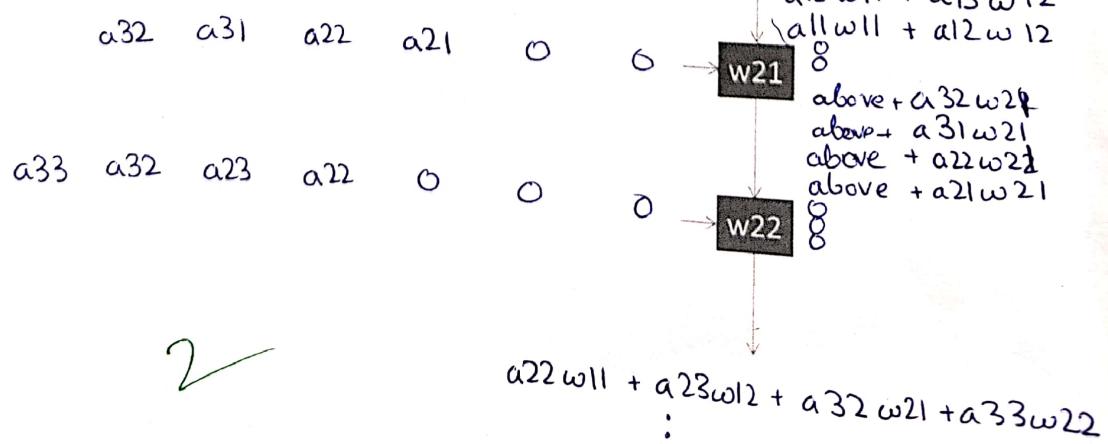
Input:

a11	a12	a13
a21	a22	a23
a31	a32	a33



Kernel:

w11	w12
w21	w22



4. An application takes 15 seconds time to execute on a multi-core machine but takes 150 seconds on a single-core machine. What could be the maximum value of the time taken by the sequential (non-parallelizable) portion of the code on the single-core machine? [2]

By Amdahl's law, (let f be sequential portion)

$$\checkmark \frac{150}{15} = \frac{f + 1-f}{f + \frac{1-f}{n}} \quad \text{as } n \rightarrow \infty \quad (\text{infinite cores})$$

$$\checkmark \Rightarrow 10 = \frac{1}{f} \Rightarrow f = 0.1 \\ = 10\% \equiv 15 \text{ seconds}$$

5. In the Eyeriss accelerator discussed in the class, data gating helps to improve: [1]

- a. Energy consumption
- b. Execution time or speed
- c. Area

6. Compress the given data as per the run-length compression scheme of Eyeriss. Show the value of term-bit as well. [2]

12, 0, 0, 0, 1, 9, 8, 0, 0, 0, 0, 0, 0, 112, 4, 0, 2, 0, 7, 0, 0, 0, 9.

The encoding format is as shown below:

Run Level Run Level Run Level Term

5b	16b	5b	16b	5b	16b	1b	
0	12	3	1	0	9	0	(next row)
0	8	6	112	0	4	0	...
1	2	1	7	3	9	1	2

Name: Pandurangi Aditya Sriram

Roll number: EE22BTECH11039

Roll number: EE22BTECH11039

Q. What is the Channel Gating network (CGNet)? Give 1-2 line reasoning.

1. Which of the following is/are true for the Channel Slicing ...
for your choices. [2]

 - a. It deploys a static pruning scheme
 - b. It provides a fixed speedup, independent of the input image
 - c. It deploys a dynamic pruning scheme
 - d. The speedup depends upon the provided input image

a. False. The pruning is dynamic based on ~~input image~~ PSUM's generated till a certain fraction of input channels. while during inference.

b. False. Speed-up is variable and depends on ~~be~~ the fraction of discriminatory information in input (lower fraction \rightarrow higher speedup).

c. True

d. True.

2. DNNExplorer uses a custom pipelined structure for initial few layers, but a generic structure for later layers. Why? [3]

 - The computation-to-communication ratio (CTR) has very high variance for earlier layers in a CNN (in general), across several orders of magnitude. This is true even for a fixed set of parameters with variable input.
 - However, for deeper layers, this variance decreases. ~~so~~
 - ~~so~~ Having a general structure throughout makes it ~~so~~ so that there is a potentially huge wastage of resources of accelerator (or) higher latency (in case of high communication), while specific structure does not allow for usage for various kinds of networks (generality)

A balance is struck by using a pipelined and generic structure each having roughly equal latency.

low variance CTC
thus can be generalized
well

Max CTC / min CTC
ratio is small

3. We observed that a few of the compressed networks are able to provide higher accuracy than original uncompressed networks. What makes this possible? [2]

In networks such as mobile net, accuracy remains similar even after a large amount of compression.

- ④ Overfitting may reduce for compressed networks with less parameters.

Sometimes, the depth-wise

For eg. separating depth-wise and point-wise convolution leads to a similar structure of different kernel channels, which may resolve input features better, as less parameters are easier to train.

- ⑤

~~Training also +~~

|

4. Explain the significance of parameters $S_{1 \times 1}$, $e_{3 \times 3}$, and $e_{1 \times 1}$ in SqueezeNet. Is there any relationship that must hold between these parameters for an effective compression? [3]

$S_{1 \times 1}$ indicates number of 1×1 squeeze kernels in a fire module.

$e_{3 \times 3}$ indicates " " 3×3 expand " " "

$e_{1 \times 1}$ " " " 1×1 " " "

- To achieve effective accuracy, the model must compress data more aggressively at later point of the network, so that the network can gather most of the useful information ~~early~~ by then.

✓

Thus, the fraction of squeeze kernels must increase for deeper layers, and expand ker and should be low for earlier layers.

~~3x3 kernels must also be~~

6.5

Hardware Architecture for Deep Learning - CS6490. Spring 2024-25.

Dept. of CSE, IIT Hyderabad

Total marks: 10

Quiz-5

Time: 12 minutes

Name: Pandragi Adiba Sriram

Roll number: EE 22 B TECH 11039

1. Which of these choices represent a tiny device considered for the MCUNet work? [1]
 - a. SRAM: 1GB, Flash: 4GB
 - b. SRAM: 256KB, Flash: 1MB
 - c. SRAM: 320KB, Flash: 100MB
 - d. SRAM: 256KB, Flash: 2MB0.5

2. Assume a hypothetical tiny device with infinite sized Flash (regular sized SRAM). In what aspects can such a device make inferencing better than a real tiny device? What limitations of existing tiny devices cannot be resolved by such a hypothetical device. [4]

Flash is a read-only memory for practical purposes on accelerators.

Advantages of inf flash :

- ① ~~All the~~ The complete code can be converted to instructions & stored. All loops can be unrolled. This significantly reduces loop overheads (branch)
- ② ~~All~~ All the weights can be stored without needing to use ~~other~~ & biases

Sparseification / quantization techniques. Thus, accuracy can be maximized.

3

Limitations :

- ① Latency can not be reduced as the ~~SRAM~~ SRAM & processing units are finite, and operations can only be reordered but not parallelized further.
- ② ~~The~~ computations
- ③ In case of large input/output activation maps, ~~the~~ tiling can not be avoided if SRAM is not large enough. This could cause computation overhead.

3. What makes training on tiny devices more challenging than inference? List (1-2 lines each) any two ideas used in MCUNetV3 to enable training on tiny devices. [3]

Training requires calculation of gradients in a directed acyclic graph. For deeper layers, more gradients need to be calculated & stored in SRAM which makes training challenging. Also, there is a large loss of accuracy when quantizing from fp32 to int8 as the calculations are different.

- ① MCUNetv3 uses ~~fp32~~ scaling of gradients by setting their max/min value to, say -128 to 127, before quantizing. This scaling factor for weights & bias is carried over for calculation & reduces inaccuracy in ~~quantized~~ Conv operation after quantization.
- ② During backprop, sparse weight updates are used to update only the most ~~com~~ important weights (chosen through a heuristic search) & biases (last K layers).

2.5

4. Mention (1-2 lines each) any three approaches used in MCUNet to reduce activation memory requirement for CNNs. [2]

v2

- ① Tiling of inputs: ~~across~~

Inputs are ~~seen~~

The first few layers of ~~input~~ CNNs usually take up the most SRAM memory as the image is of high resolution. The no of input pixels reduces with layers.

For such cases, the ~~smaller~~ MCUNetv2 tiles input while

- 0.5
② reducing ~~is~~ SRAM memory & minimizing computation overhead.

↓
into multiple cycles

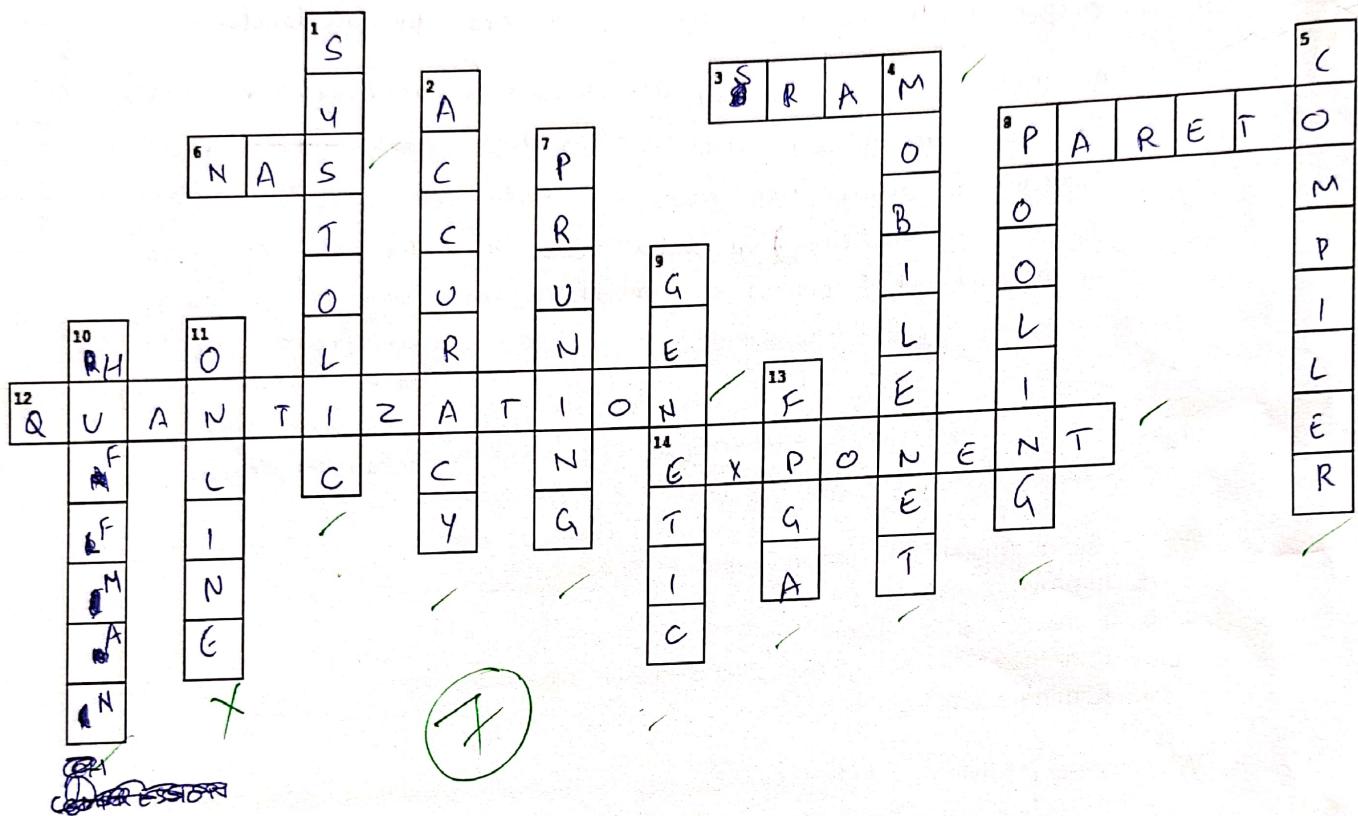
③

14

Hardware Architecture for Deep Learning - CS6490. Spring 2024-25.
Dept. of CSE, IIT Hyderabad

Quiz-6**Time: 12 minutes****Total marks: 15****Name:** P Aditya Srivastava**Roll number:** EE22BTECH11039

1. Solve the given crossword puzzle. [8]

**Across**

- 3. The place to store activations in a tiny device
- 6. Process of identify suitable CNN architecture
- 8. The non-dominated points during DSE
- 12. Using INT8 instead of FP32 representation
- 14. Determines the range of a floating point number

Down

- 1. An array structure used in most CNN accelerators
- 2. A common objective metric used in NAS
- 3. A CNN which proposed depthwise and pointwise convolutions
- 4. Helps map a CNN graph efficiently to given hardware
- 5. Removing unimportant weights from a CNN model
- 6. The layer which helps reduce the input dimensions for next layer
- 7. A commonly used bio-inspired algorithm for DSE
- 8. A lossless compression scheme used in Eyeriss Deep compression
- 9. Accelerators help to reduce this metric
- 10. A device to implement custom digital logic

2. Briefly explain the importance of using a dual-port SRAM in accelerators. [3]

Accelerators' PEs need to access SRAM to load or read from, as a global memory with the accelerator. It acts as a double buffer b/w accelerator & main (DRAM) memory.

Usually, the weights, if and of maps are loaded/saved to SRAM memory. One port is required for input read and output write to main memory & one for accelerator.

3

As instructions (high-level) are issued to accelerator (or) in case accelerator card is inserted into ~~memory~~ a system & reads from memory, data gets stored in SRAM (to reduce latency) as data reuse is extremely common in neural nets & CNNs) & simultaneously may be written to

3. Systolic arrays are typically implemented on: [2] or read from the PEs (or their local memory).
- a. FPGA
 - b. GPU
 - c. ASIC
 - d. CPU with Vector instructions

It acts like a larger "cache".

2

4. What is the purpose of mutation in Genetic Algorithms? [1]

- a. Improve crossover
- b. Increase the population
- c. Introduce new characteristics not present in initial population
- d. Eliminate constraints

1

5. Which of the following are examples of bio-inspired heuristic optimization algorithms? [1]

- a. Particle Swarm Optimization
- b. Genetic Algorithm
- c. Gradient Descent
- d. Ant Colony Optimization

1