

Space

Eq

Compute

Complexity:

Space required to store the model's weights.

+
Space required to store the activations of 1 forward pass.

$$y = \underline{w}x + b$$

operations.

(a) \Rightarrow addition
 \downarrow
 $wx + b$
(b) \Rightarrow assignment
 \hookrightarrow
 $wx + b$ into y .
(c) \Rightarrow multiplication
 \hookrightarrow wx .
MAC



minimum (bs=1)

SC \Rightarrow how much memory is required to do 1 forward of an input?



Space for parameters

+

Space for activations

Space

$$(l_1 + l_2)$$

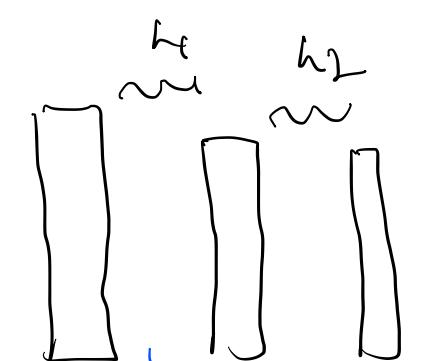
\Downarrow

$$(l_1 + l_2) + \underline{(x_0 + x_1 + x_{op})}$$

\Downarrow

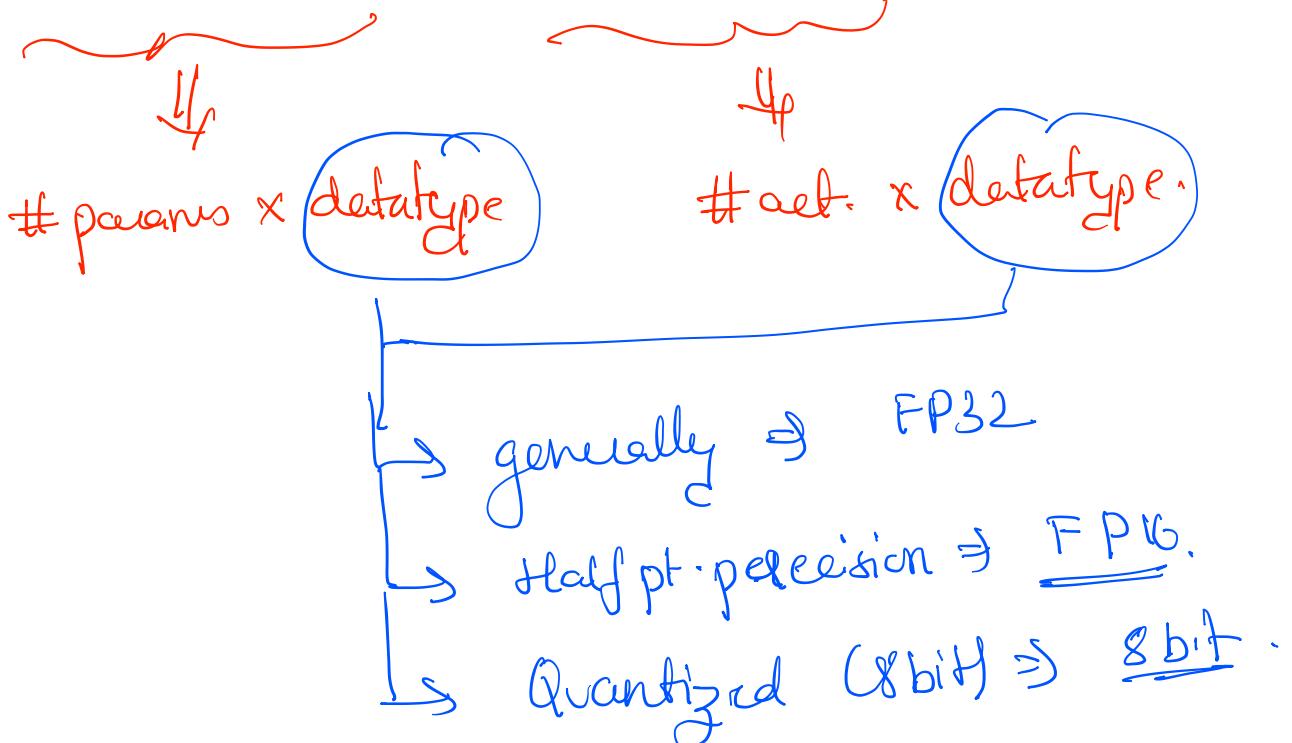
parameters

activations.
activations.

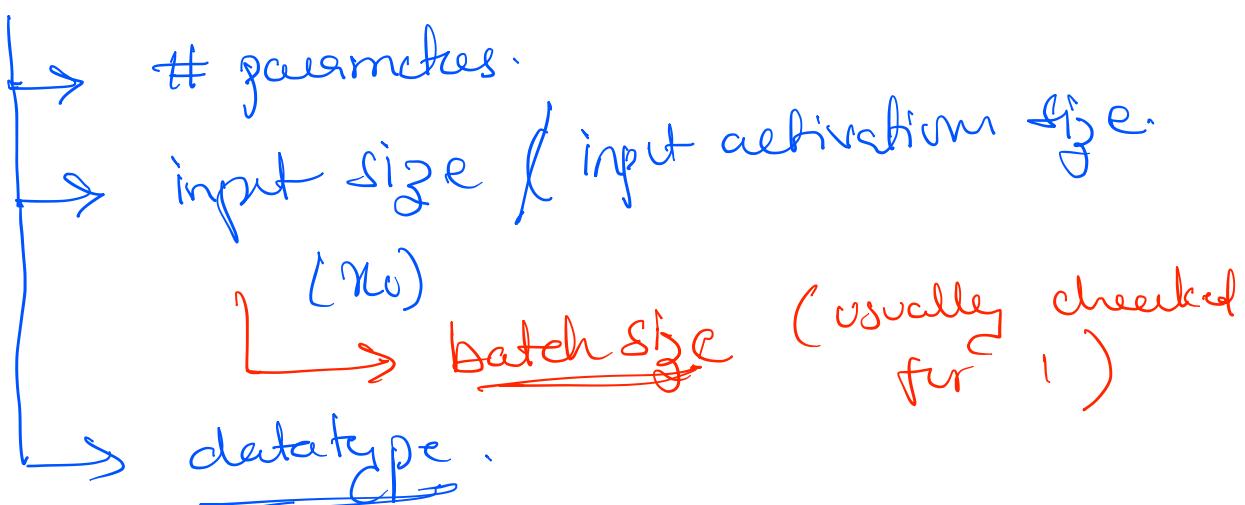


activations

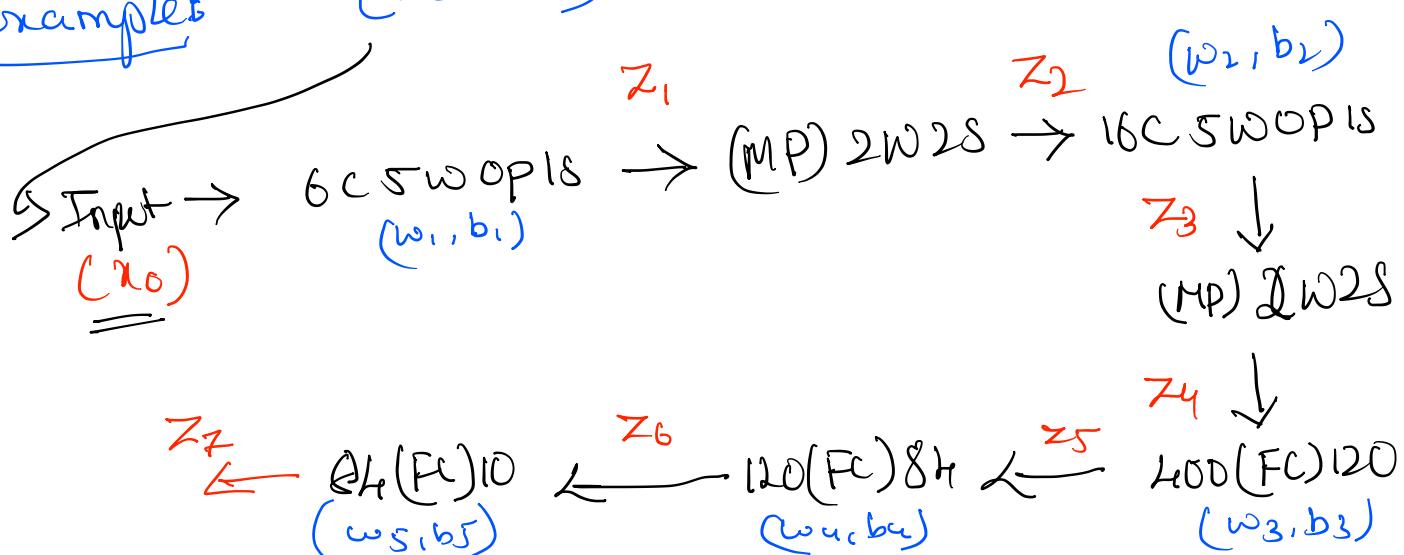
weights.
(including bias).



Space complexity: is affected by -



Example: (lenet5)



Consider: If $I \Rightarrow (x_0)$ \Downarrow $I \in 1 \times 32 \times 32$
 represented as $\underline{x_0}$ (writing as x_0 for simplicity)

Activations

$\text{Size}(x_0) = 1 \times 32 \times 32 = 1,024$

$\text{Size}(z_1) = 5 \times 28 \times 28$
 $= \frac{32 - 5 + 2(0)}{4} + 1 = 28$

$\text{Size}(z_2) = 6 \times 14 \times 14 = 1,176.$

$\text{Size}(z_3) = 1600$

$\text{Size}(z_4) = 400$

$\text{Size}(z_5) = 120$

$\text{Size}(z_6) = 84$

$\text{Size}(z_7) = 10$

$\text{Total}_{\text{act}} = 8,094.$

$\hookrightarrow \text{Size} = (8,094) \times \frac{\text{datatype}}{\text{C}}$
 usually (\hookrightarrow FP32)

Parameters

$\text{Size}(w_1, b_1) =$

derivation. \hookrightarrow CNN formulae.

$I \in 1 \times M \times N$

$K \in 2 \times (1 \times \text{width})$

$w \in 2 \times M_1 \times N_1$

$M_1 = \left\lceil \frac{M - w + 2p_w}{sw} + 1 \right\rceil$

$N_1 = \left\lceil \frac{N - h + 2p_h}{sh} + 1 \right\rceil$

Where is it coming from?

$\# \text{Kernels in } K$

$\Rightarrow \# \text{output channels} = \# \text{independent kernels in } K$

$\hookrightarrow \text{Size of each such kernel?}$

$$\left[(\# \text{ o/p channels}) \times (\# \text{ i/p channels}) \right] \times w \times h$$

$$(6 \times 1 \times 5 \times 5)$$

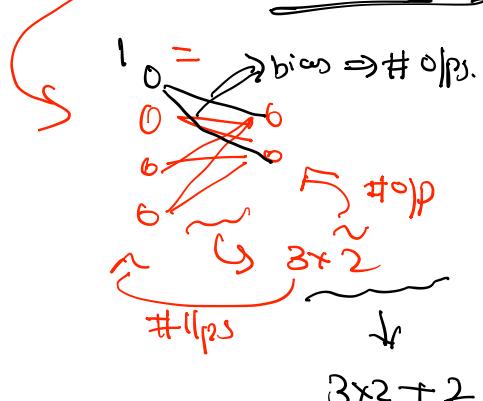
$$\text{Size } (\omega_1) = 150.$$

$$\begin{aligned} \text{Size } (\omega_1) &= \# \text{ o/p channels} \\ &= 6 \end{aligned}$$

$$\begin{aligned} \text{Size } (\omega_1, b_1) &= 150 + 6 \\ &= 156 // \end{aligned}$$

$$\begin{aligned} \text{Size } (\omega_2, b_2) &= 16 \times 6 \times 5 \times 5 + 16 \\ &= 2416 \end{aligned}$$

$$48,120 = \text{Size } (\omega_3, b_3) = 400 \times 120 + 120 \Rightarrow (\underbrace{\text{i/p} \times \text{o/p} + \text{o/p}}_{\# \text{ o/p}})$$



$$\begin{aligned} \text{Size } (\omega_4, b_4) &= 80 \times 84 + 84 \\ &= 10,164 \end{aligned}$$

$$\begin{aligned} \text{Size } (\omega_5, b_5) &= 84 \times 10 + 10 \\ &= 850 \end{aligned}$$

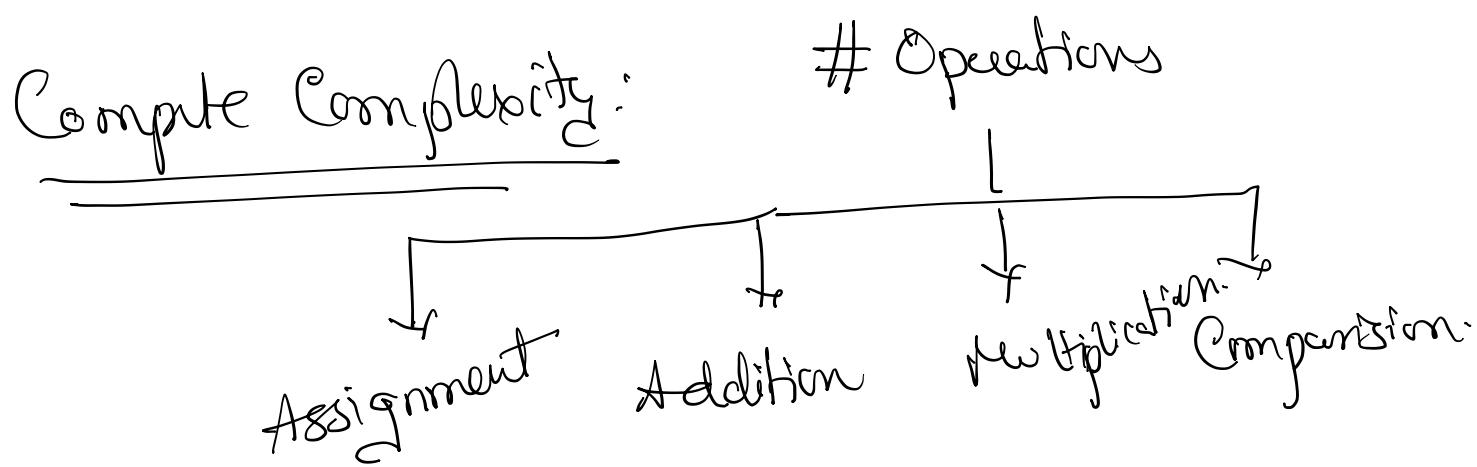
$$\text{Total params} = 156 + 2416 + 48,120 + 10,164 + 850$$

$$\text{size} = \frac{\text{Total params}}{\text{datatype}}$$

usually FP32

Total Space: $(\text{Total}_{\text{act}} + \text{Total}_{\text{params}}) \times \text{datatype}$

$$= (8,094 + 61,706) \times (\text{dt})$$
$$\approx \underline{\underline{69k}} \times (\text{dt})$$



~~#~~ for a perceptron:

$$y = w^T x + b$$

Op activation: $(|x|)$

$\uparrow (|x|)$ $\uparrow (|x|)$ $\uparrow (|x|)$

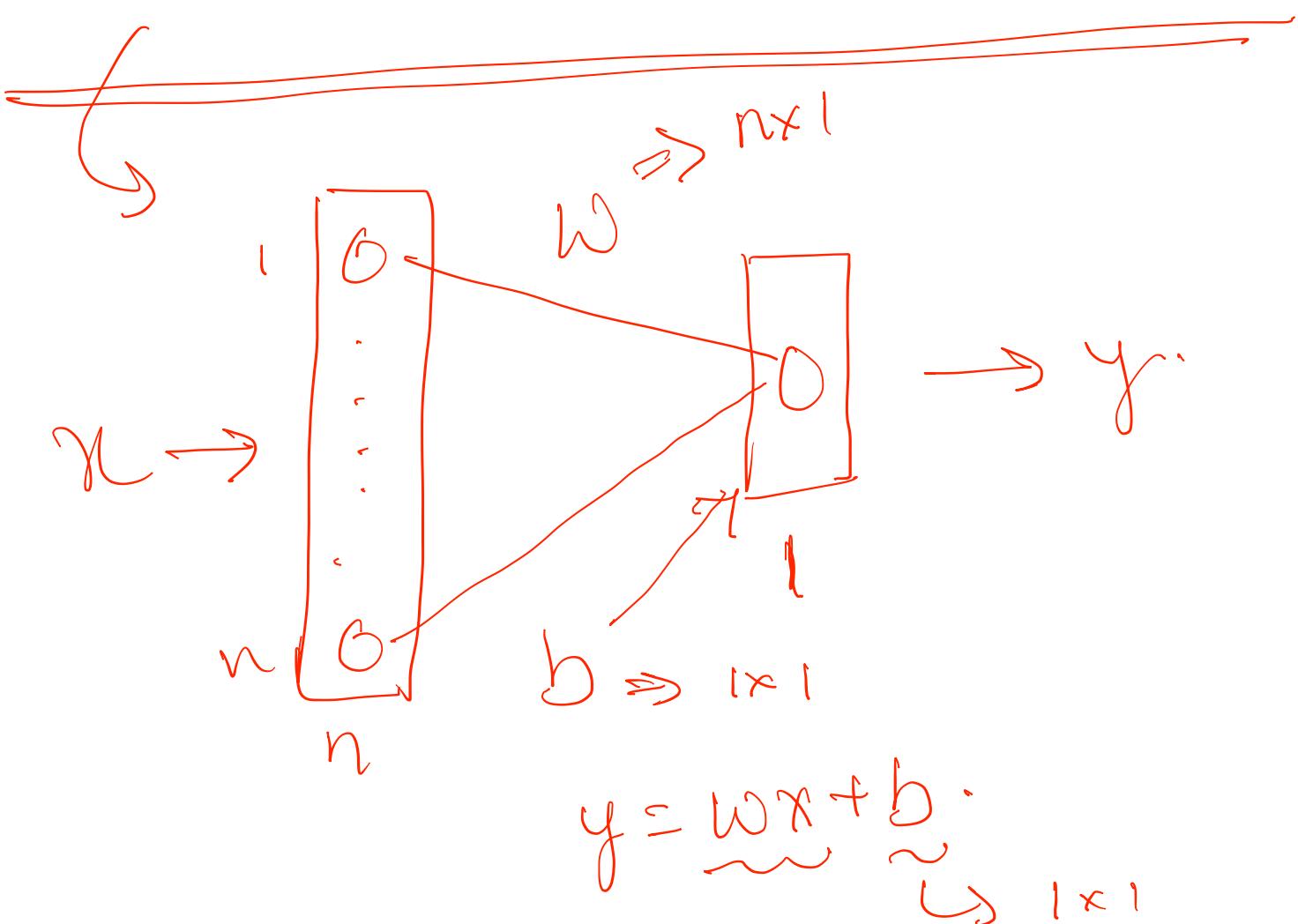
\rightarrow if activation: $(|x|)$

$\left\{ \begin{array}{l} \text{Multiplication: } (n) \\ \text{Addition: } (A) \\ \text{Assignment: } (implied) \end{array} \right.$

non-pipelined
 computation
 ↪
 no-parallel
 processing.

1 op activation
 \hookrightarrow
 $G(M + IA) \Rightarrow \underline{\text{2 operations}}$

\hookrightarrow per op. activation
 depends on size
 $\text{of input} \Rightarrow 0$



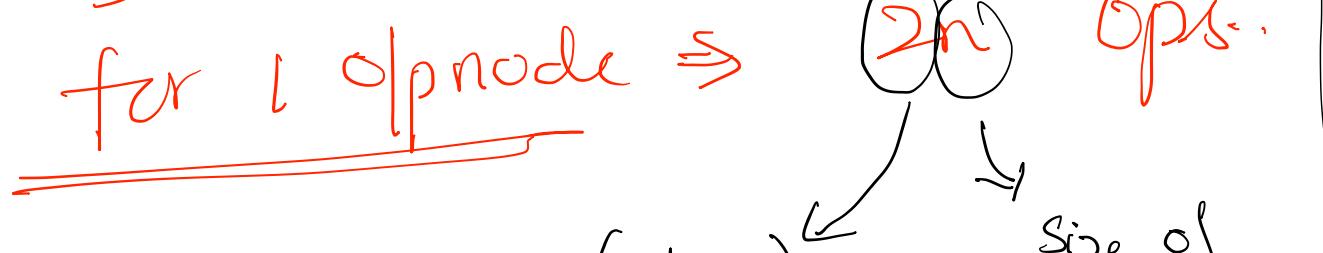
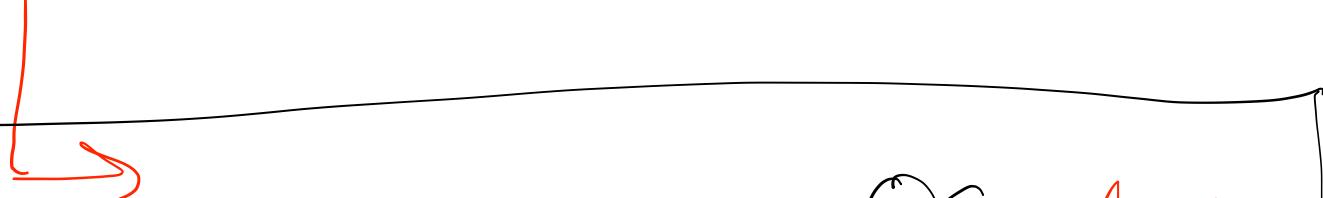
$$wx \Rightarrow \sum_{i=1}^n w_i x_i$$

$\hookrightarrow n \rightarrow \text{multiplications}$

$(n-1) \rightarrow \text{additions.}$

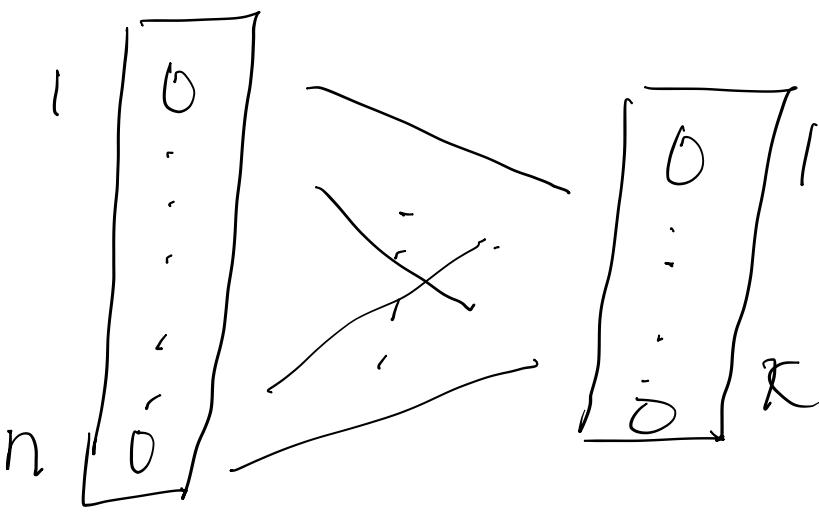
$$\underline{t_b} \Rightarrow l \rightarrow \text{addition}$$

total operations $\Rightarrow (nM + nA)$
 $= 2n \text{ operations}$



$(M + A)$ $\xrightarrow{\quad}$ $2n$ Ops.
 Size of input.

(non pipelined)



General
NN

$$\# \text{ operations} = k(2n)$$

↓
 # of nodes ↓
 # of pip nodes
 (M+A)

non-pipelined \Rightarrow M & A happen sequentially.

Pipelined \Rightarrow parallel compute.
 (FMA block is used) \hookrightarrow M & A happen parallelly

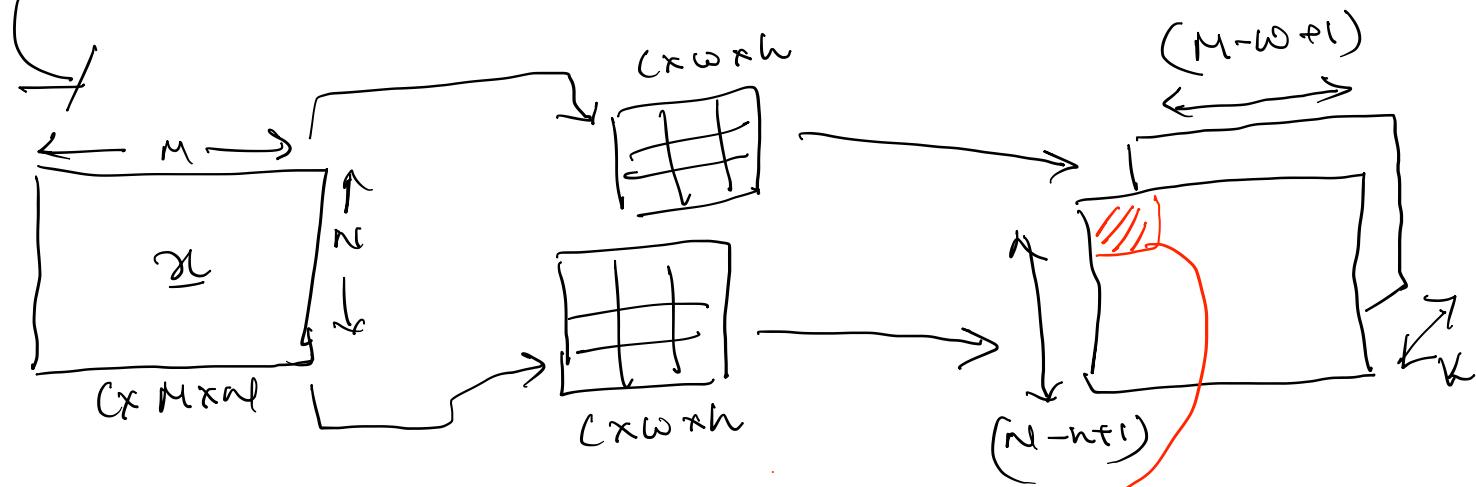
(pipelined) \hookrightarrow

$$\# \text{ ops} = k(\tilde{n} + 1)$$

Compute Complexity of CNNs:

Consider:

→ An Input, $x \in \mathbb{R}^{C \times M \times N}$
 $(s=1, p=0)$ → A conv. kernel, $w \in \mathbb{R}^{K \times C \times W \times H}$
 → o/p activation, $z \in \mathbb{R}^{K \times (M-W+1) \times (N-H+1)}$



for 1 op node $\Rightarrow z_{k,x,y} = w_{k,0} + \sum_{i=1}^w \sum_{j=1}^h w_{k,i,j} x_{(x+i, y+j)}$

$\boxed{2chw \text{ operations}} \Leftrightarrow 1A + chwMs + chw - 1 As$

$\hookrightarrow \# \text{op nodes} = K \times (M-W+1) \times (N-H+1) \quad (\because s=0, p=1)$

$K \times \left(\frac{M-W+2P_w}{S_w} + 1 \right) \left(\frac{N-H+2P_h}{S_h} + 1 \right)$ else

→ Total # operations = $\left[k \times \left(\frac{M-w+2Pw}{sw} + 1 \right) \left(\frac{N-h+2Ph}{sn} + 1 \right) \right] [2chw]$

↓
slp nodes

+
ops per op node.

→ Total ops (pipelined) = $\left[k \times \left(\frac{M-w+2Pw}{sw} + 1 \right) \left(\frac{N-h+2Ph}{sn} + 1 \right) \right] (chw+1)$

↓
dp nodes

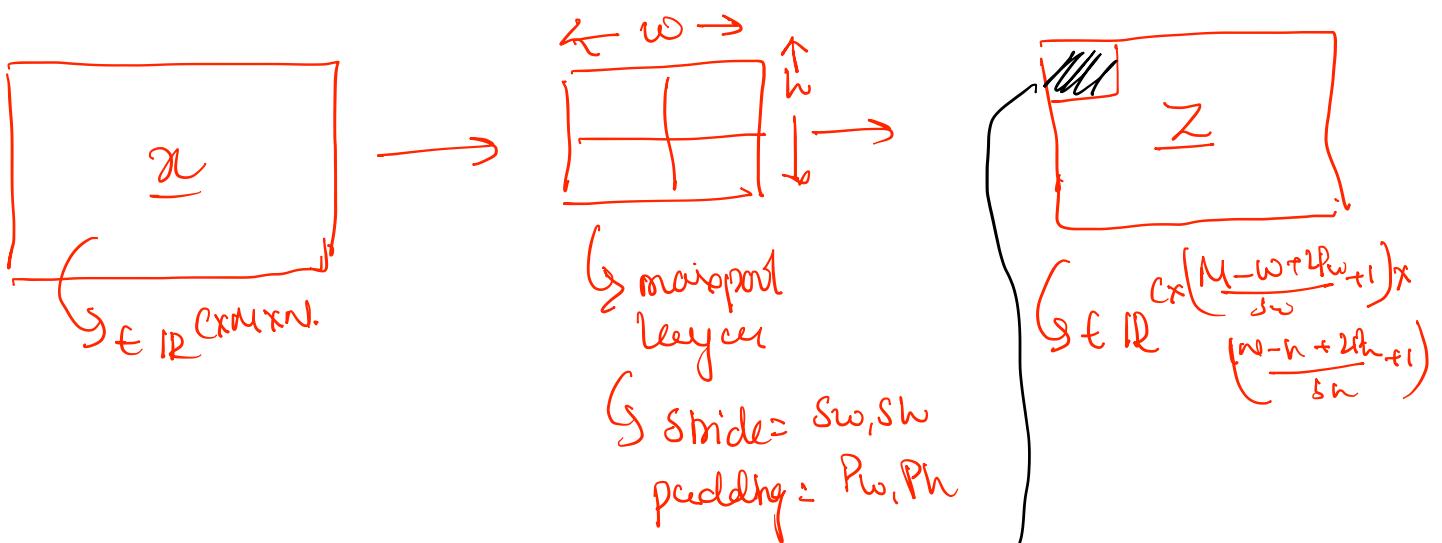
↓
ops per op node

L1 & L2 cache
parallel using FMA
block

Fused Multiply and Add H/w Acceler.

→ Supports parallel addition & multiplication operations.

Compute Complexity of Maxpool layer:



per 1 op node: $(hw - 1)$ comparison operations.

Total ops \Rightarrow
$$\left[C \times \left(\frac{H - w + 2p_w + 1}{s_w} \right) \left(\frac{W - w + 2p_h + 1}{s_h} \right) \right] [hw - 1]$$

↓
#Ops per node.

#Op nodes

Speed :- # operations (Floating point operations per second)
 (in order of GFLOPs) \propto FLOPs (based on processor's clock)

Latency : Time taken to process an input (image).

$$\text{Latency} = \frac{\# \text{ Total operations}}{\text{Speed.}} \quad (\text{units} = \text{seconds})$$

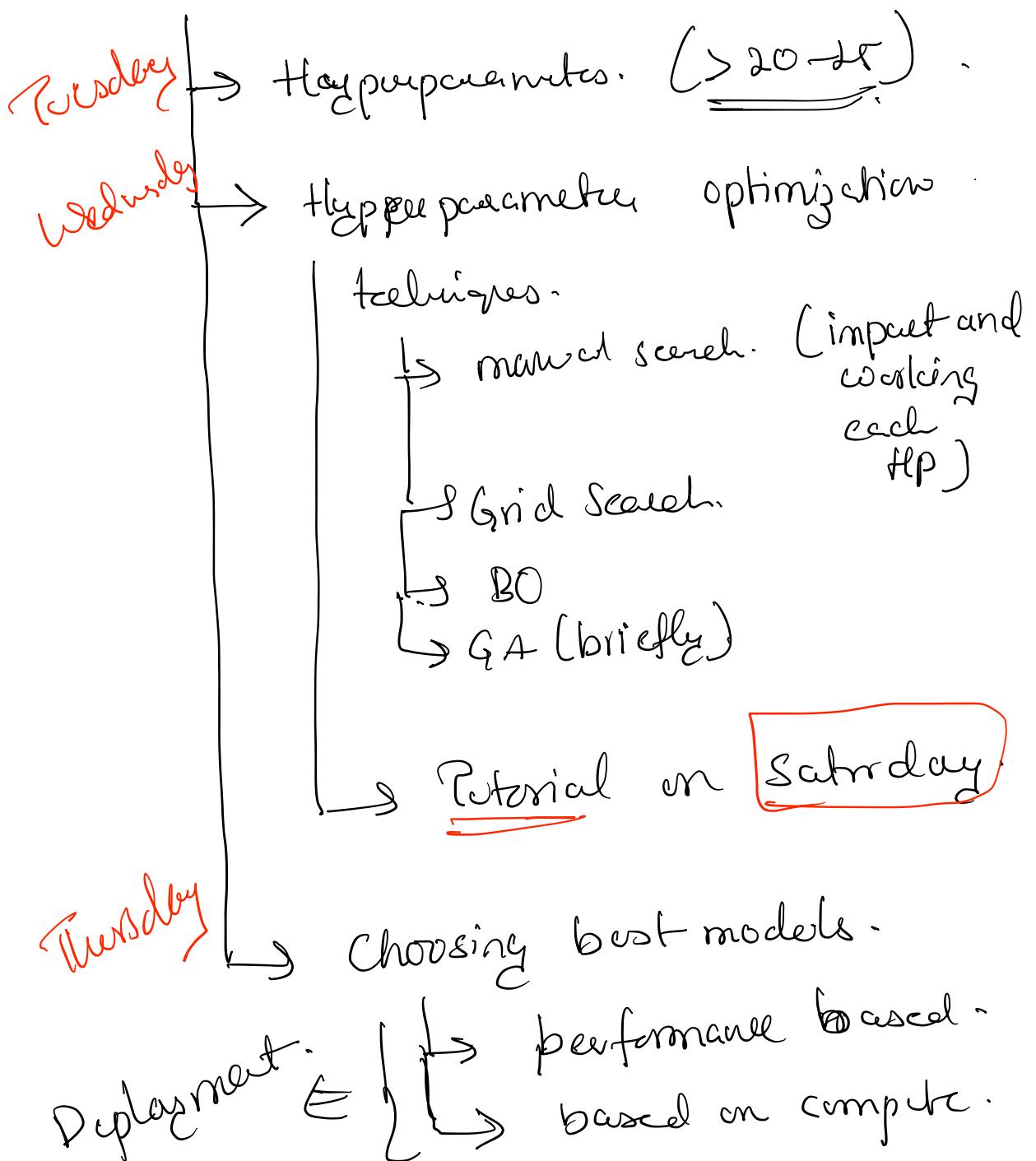
Throughput : # Images processed / unit time.

$$T = \frac{\text{Speed.}}{\# \text{ Total Operations}} \quad (\text{units} \Rightarrow \frac{\text{images}}{\text{second}})$$

Assuming all images are processed sequentially.

Next week \Rightarrow Final week: (Tuesday - Saturday).

Advanced Concepts:



Topic of your choice:

FRIDAY