

# Fault Tolerant CNN

# Error Prone CNN Application

- Recent literature indicates that soft errors are inevitable in modern systems, from edge computing devices to supercomputers
  - high-energy cosmic radiation
  - Aging and wear of devices
- CNN inference applications often call for power-efficient and cost-efficient machine learning accelerators, which may adopt overclocking with voltage undervolting, incurring more soft errors than common hardware incurs.

- Resilient convolutional neural networks are essential for guaranteeing the correctness of inference applications.
- A single bit flip happened during CNN image classification could result in as much as 40% and 70% SDC rate in datapath and memory, respectively.
- Such high SDC rates would downgrade the CNN prediction accuracy dramatically.

### **[SDC Rate: Silent Data Corruption (SDC) Rate]**

- Existing resilient solutions (ECC) are insufficient for protecting CNN inference applications against multiple soft errors.

# CNN Basics

$$\mathbf{O}[n][m][x][y] = \mathbf{B}[m] + \sum_{k=0}^{Ch-1} \sum_{i=0}^{R-1} \sum_{j=0}^{R-1} \mathbf{D}[n][k][Ux+i][Uy+j] \times \mathbf{W}[m][k][i][j]$$

$$0 \leq n < N, 0 \leq m < M, 0 \leq x, y < E, E = \frac{H-R+U}{U}$$

The convolution operation involves two significant inputs: the feature map (fmap)  $\mathbf{D}$ ,  $\mathbf{D} \in \mathbb{R}^{N \times Ch \times H \times H}$ , and the convolutional kernels  $\mathbf{W}$ ,  $\mathbf{W} \in \mathbb{R}^{M \times Ch \times R \times R}$ .

$\mathbf{B}$ , is applied to the output after convolution, and the final result is denoted as  $\mathbf{O}$ ,  $\mathbf{O} \in \mathbb{R}^{N \times M \times E \times E}$ .

Notation	Description
$\mathbf{D}$	Feature map, dimension is $4D$
$\mathbf{W}$	Kernels, also called filters, dimension is $4D$
$\mathbf{O}$	Output, dimension is $4D$
$\mathbf{B}$	Bias, dimension is $1D$
$\mathbf{C}$	Checksums
$\mathbf{S}$	Block summations of $\mathbf{O}$ , corresponding to checksums
$\otimes$	Convolution operation
$N$	First dimension of $\mathbf{D}$ and $\mathbf{O}$
$M$	First dimension of $\mathbf{W}$ and second dimension of $\mathbf{O}$
$Ch$	Second dimension of $\mathbf{D}$ and $\mathbf{W}$ , also called channels
$H$	Third and fourth dimension of $\mathbf{D}$
$R$	Third and fourth dimension of $\mathbf{W}$
$E$	Third and fourth dimension of $\mathbf{O}$
$U$	Stride size

# CONV Layer Tensor

Output fmaps (O)

Computation

Input fmaps (I)

Biases (B)

Filter weights (W)

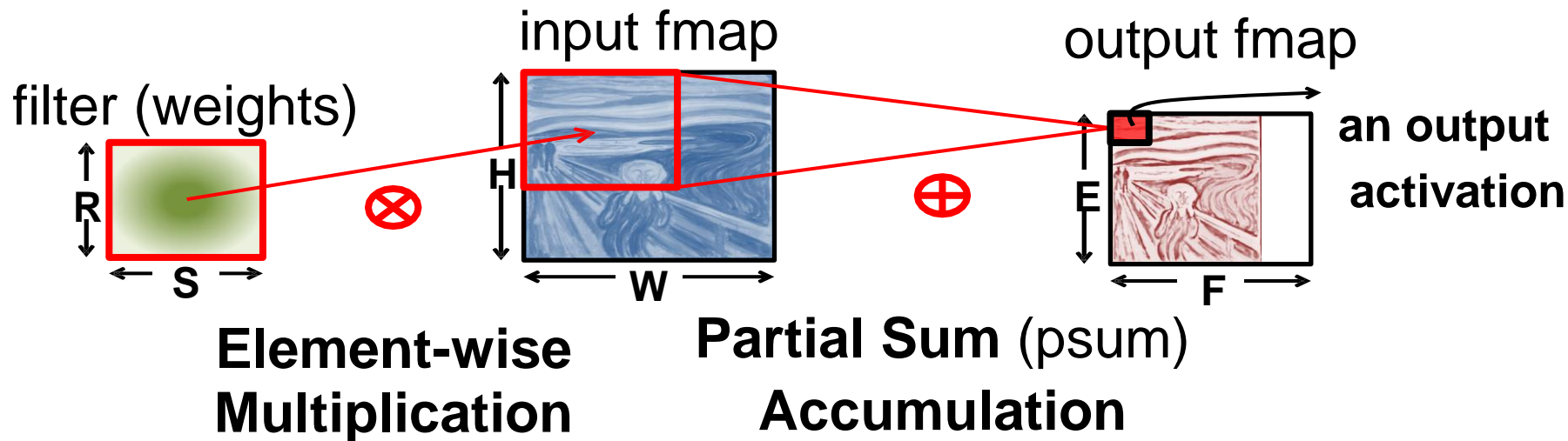
$$\underline{\mathbf{O}[n][m][x][y]} = \text{Activation}(\underline{\mathbf{B}[m]} + \sum_{i=0}^{R-1} \sum_{j=0}^{S-1} \sum_{k=0}^{C-1} \underline{\mathbf{I}[n][k][Ux+i][Uy+j]} \times \underline{\mathbf{W}[m][k][i][j]}),$$

$$0 \leq n < N, 0 \leq m < M, 0 \leq y < E, 0 \leq x < F,$$

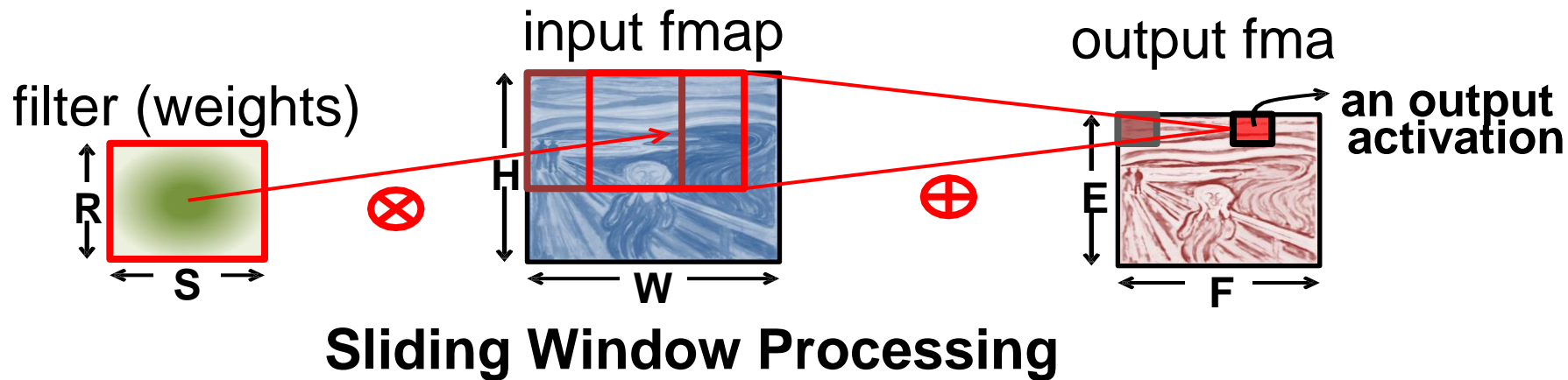
$$E = (H - R + U)/U, F = (W - S + U)/U.$$

Shape Parameter	Description
$N$	fmap batch size
$M$	# of filters / # of output fmap channels
$C$	# of input fmap/filter channels
$H/W$	input fmap height/width
$R/S$	filter height/width
$E/F$	output fmap height/width
$U$	convolution stride

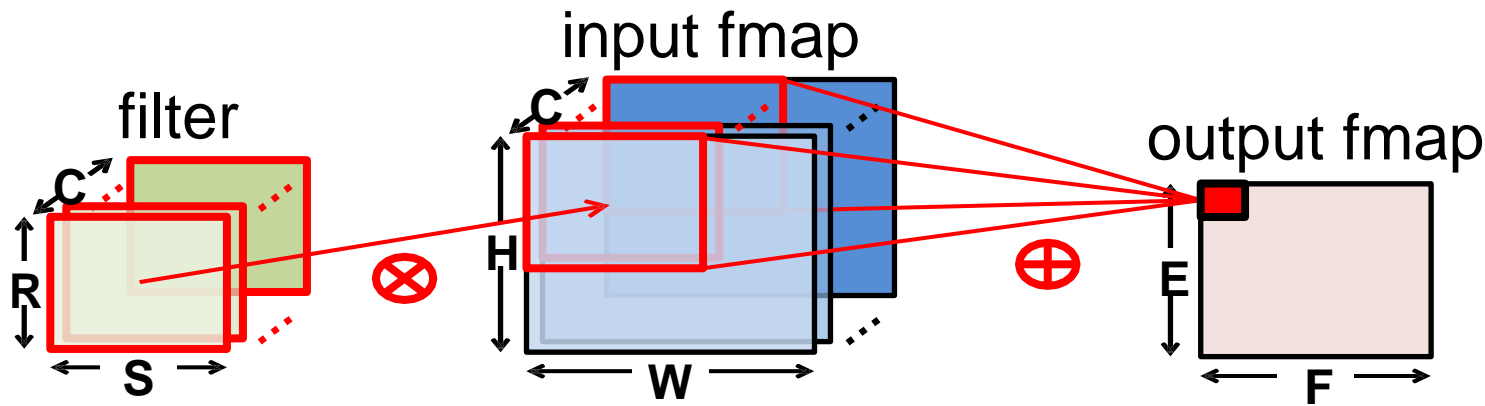
# Convolution (CONV) Layer



# Convolution (CONV) Layer



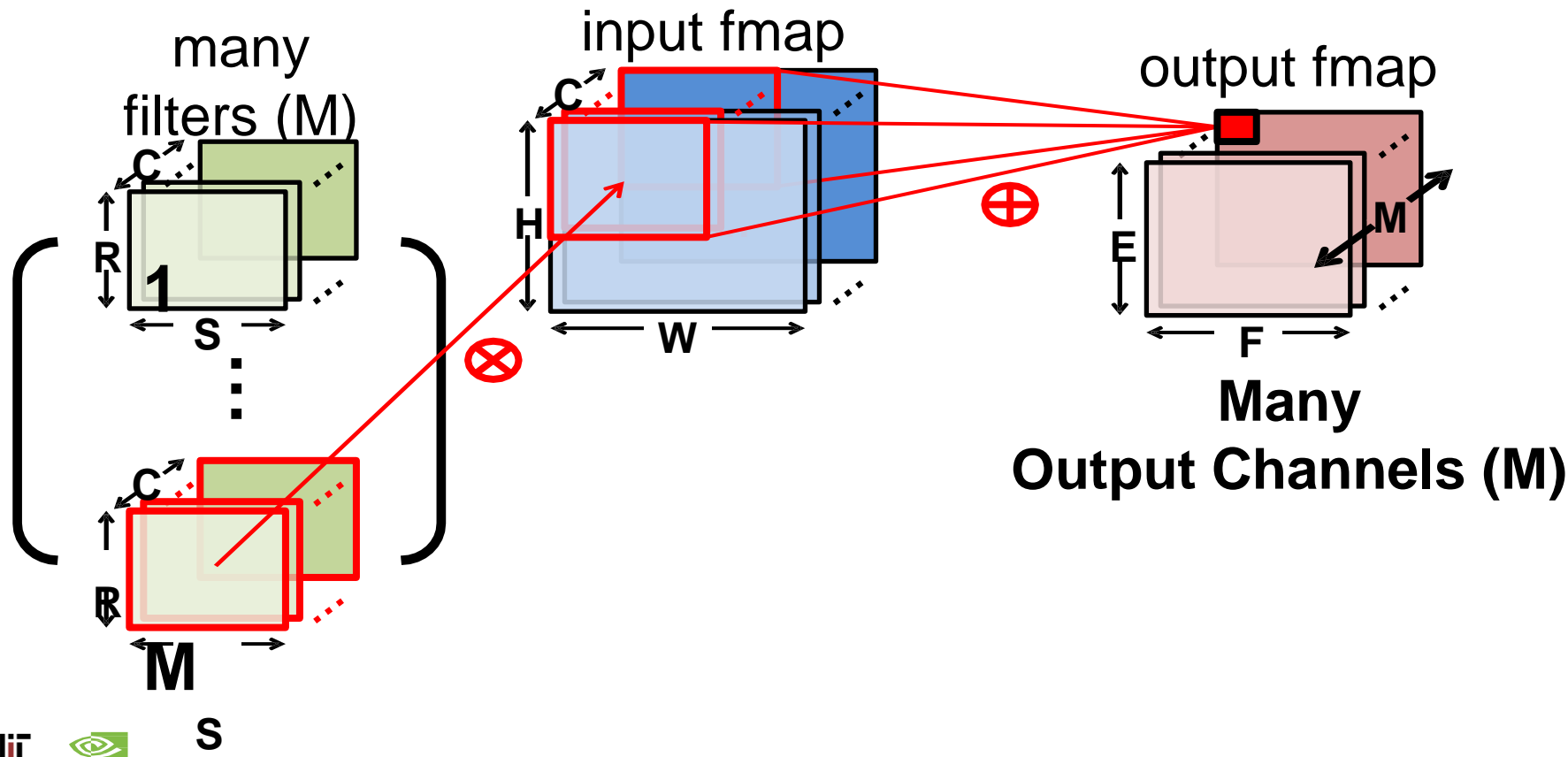
# Convolution (CONV) Layer



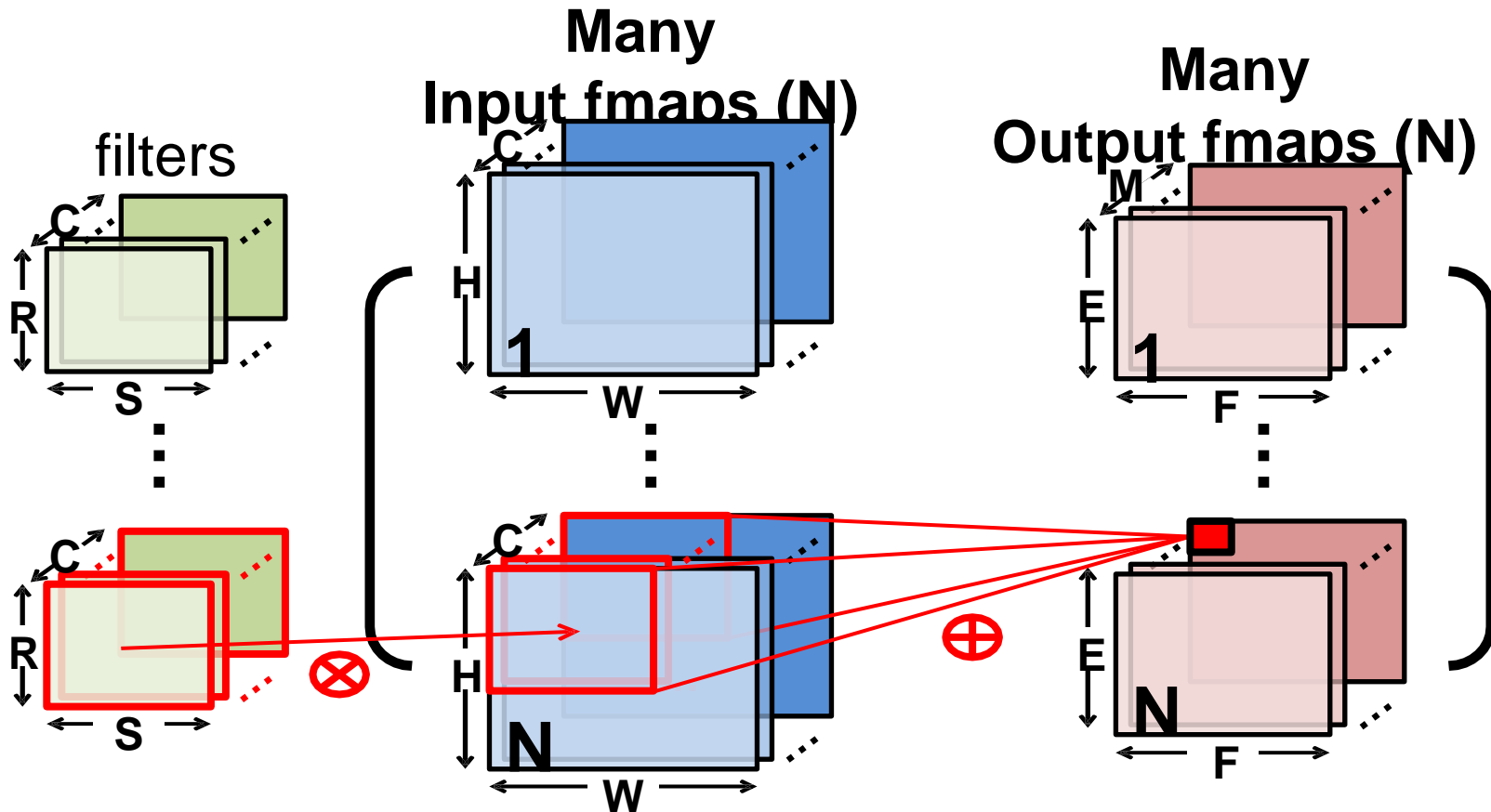
**Many Input Channels ( $C$ )**



# Convolution (CONV) Layer



# Convolution (CONV) Layer



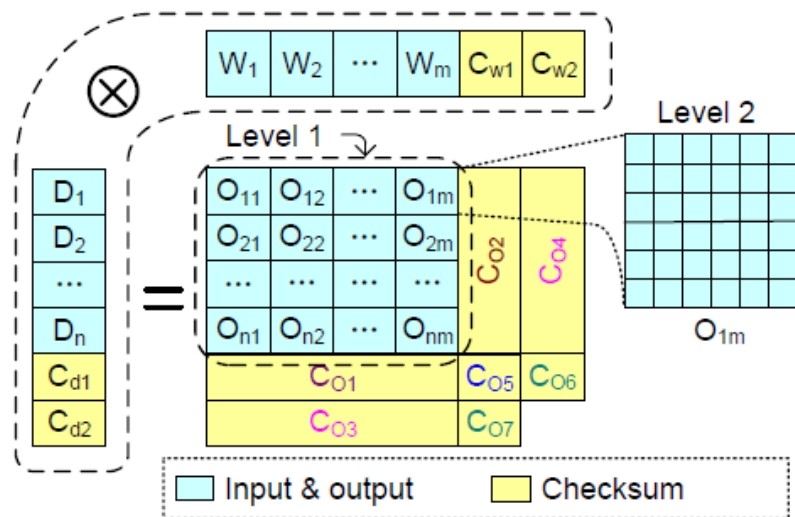
# CNN Decoder Ring

---

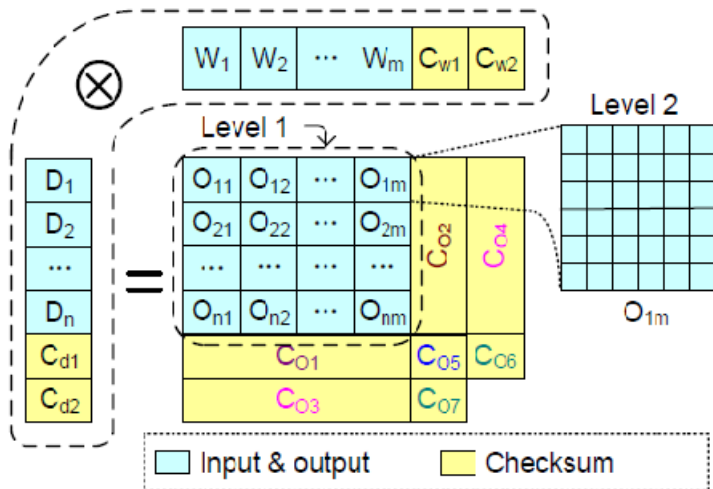
- N – Number of **input fmaps/output fmaps** (batch size)
- C – Number of 2-D **input fmaps /filters** (channels)
- H – Height of **input fmap** (activations)
- W – Width of **input fmap** (activations)
- R – Height of 2-D **filter** (weights)
- S – Width of 2-D **filter** (weights)
- M – Number of 2-D **output fmaps** (channels)
- E – Height of **output fmap** (activations)
- F – Width of **output fmap** (activations)

# Checksum Support

Notation	Description
$D$	Feature map, dimension is $4D$
$W$	Kernels, also called filters, dimension is $4D$
$O$	Output, dimension is $4D$
$B$	Bias, dimension is $1D$
$C$	Checksums
$S$	Block summations of $O$ , corresponding to checksums
$\otimes$	Convolution operation
$N$	First dimension of $D$ and $O$
$M$	First dimension of $W$ and second dimension of $O$
$Ch$	Second dimension of $D$ and $W$ , also called channels
$H$	Third and fourth dimension of $D$
$R$	Third and fourth dimension of $W$
$E$	Third and fourth dimension of $O$
$U$	Stride size



# CNN Checksums



$$C_{d1} = \sum_{n=0}^{N-1} D_n$$

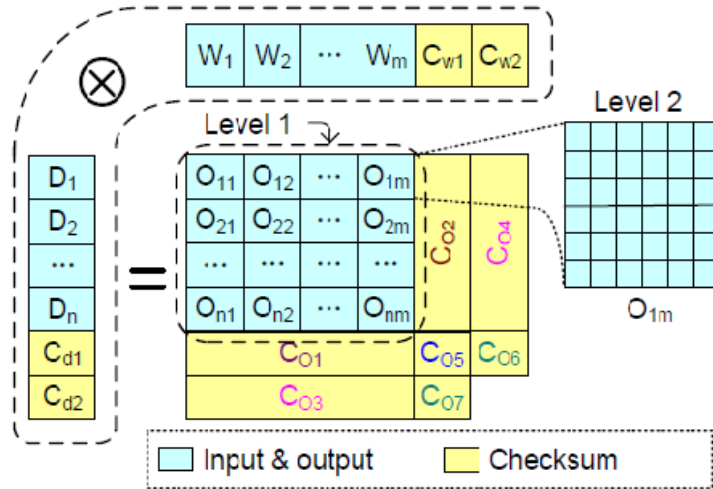
$$C_{d2} = \sum_{n=0}^{N-1} n D_n$$

$$C_{w1} = \sum_{m=0}^{M-1} W_m$$

$$C_{w2} = \sum_{m=0}^{M-1} m W_m$$

- $C_{d1}$  and  $C_{w1}$  are calculated before the convolution operation,
- Any memory error striking  $D$  or  $W$  during the convolution would not affect  $C_{d1}$  or  $C_{w1}$ .

# CNN Checksums

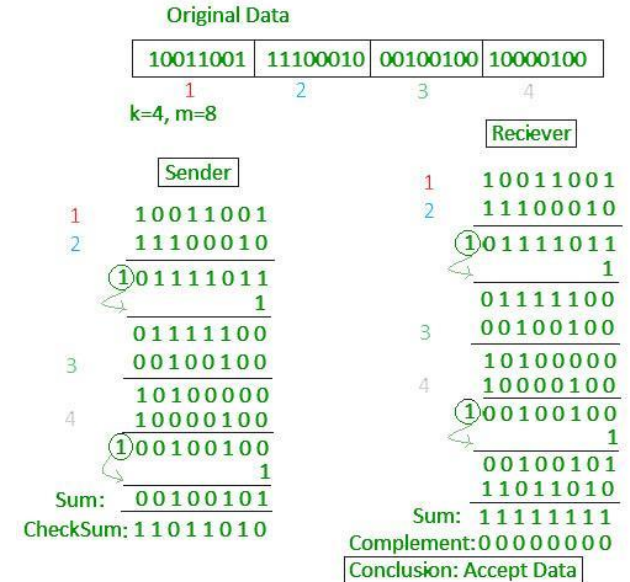


$$\begin{aligned} C_{o1} &= C_{d1} \otimes W \\ C_{o2} &= D \otimes C_{w1} \\ C_{o3} &= C_{d2} \otimes W \\ C_{o4} &= D \otimes C_{w2} \\ C_{o5} &= C_{d1} \otimes C_{w1} \\ C_{o6} &= C_{d1} \otimes C_{w2} \\ C_{o7} &= C_{d2} \otimes C_{w1} \end{aligned}$$

- The output  $O$  is represented in the form of blocks.
- Elements inside the same block are independent with respect to checksums.
- perform the checksum comparison independently for each element across blocks.
- Multiple soft errors in the same block can be detected and corrected independently.

# Original Checksum

- In checksum error detection scheme, the data is divided into  $k$  segments each of  $m$  bits.
- In the sender's end the segments are added using 1's complement arithmetic to get the sum. The sum is complemented to get the checksum.
- The checksum segment is sent along with the data segments.
- At the receiver's end, all received segments are added using 1's complement arithmetic to get the sum. The sum is complemented.
- If the result is zero, the received data is accepted; otherwise discarded.



# Distributive property of Convolution

$$\begin{aligned}
 & \mathbf{D}_1 \otimes \mathbf{W}_1 + \mathbf{D}_2 \otimes \mathbf{W}_1 = \\
 & \sum_{k=0}^{Ch-1} \sum_{i=0}^{R-1} \sum_{j=0}^{R-1} \mathbf{D}_1[k][Ux+i][Uy+j] \times \mathbf{W}_1[k][i][j] \\
 & + \sum_{k=0}^{Ch-1} \sum_{i=0}^{R-1} \sum_{j=0}^{R-1} \mathbf{D}_2[k][Ux+i][Uy+j] \times \mathbf{W}_1[k][i][j] \\
 & = \sum_{k=0}^{Ch-1} \sum_{i=0}^{R-1} \sum_{j=0}^{R-1} (\mathbf{D}_1 + \mathbf{D}_2)[k][Ux+i][Uy+j] \times \mathbf{W}_1[k][i][j] \\
 & = (\mathbf{D}_1 + \mathbf{D}_2) \otimes \mathbf{W}_1
 \end{aligned}$$

Similarly, we can get the following equation.

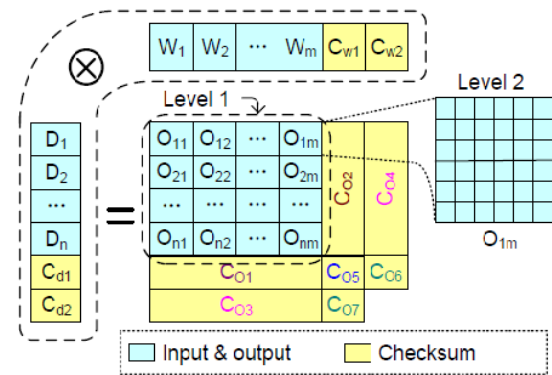
$$\mathbf{D}_1 \otimes \mathbf{W}_1 + \mathbf{D}_1 \otimes \mathbf{W}_2 = \mathbf{D}_1 \otimes (\mathbf{W}_1 + \mathbf{W}_2)$$

Notation	Description
$\mathbf{D}$	Feature map, dimension is $4D$
$\mathbf{W}$	Kernels, also called filters, dimension is $4D$
$\mathbf{O}$	Output, dimension is $4D$
$\mathbf{B}$	Bias, dimension is $1D$
$\mathbf{C}$	Checksums
$\mathbf{S}$	Block summations of $\mathbf{O}$ , corresponding to checksums
$\otimes$	Convolution operation
$N$	First dimension of $\mathbf{D}$ and $\mathbf{O}$
$M$	First dimension of $\mathbf{W}$ and second dimension of $\mathbf{O}$
$Ch$	Second dimension of $\mathbf{D}$ and $\mathbf{W}$ , also called channels
$H$	Third and fourth dimension of $\mathbf{D}$
$R$	Third and fourth dimension of $\mathbf{W}$
$E$	Third and fourth dimension of $\mathbf{O}$
$U$	Stride size



$$C_{o1}[m] = \left( \sum_{n=0}^{N-1} D_n \right) \otimes W_m = \sum_{n=0}^{N-1} (D_n \otimes W_m) = \sum_{n=0}^{N-1} O_{nm}$$

$$C_{o2}[n] = D_n \otimes \left( \sum_{m=0}^{M-1} W_m \right) = \sum_{m=0}^{M-1} (D_n \otimes W_m) = \sum_{m=0}^{M-1} O_{nm}$$



These equations show the equality between the sum of output and the output checksums. Let  $S_{o1}$  and  $S_{o2}$  be the summation of the output, where  $S_{o1}[m] = \sum_{n=0}^{N-1} O_{nm}$ ,  $S_{o2}[n] = \sum_{m=0}^{M-1} O_{nm}$ . We can compare  $C_{o1}$ ,  $C_{o2}$  with  $S_{o1}$ ,  $S_{o2}$  to detect, locate, and correct soft errors if they exist.

# Row Checksum Scheme (RC)

$$C_{o1} = C_{d1} \otimes W$$

$$C_{o2} = D \otimes C_{w1}$$

$$C_{o3} = C_{d2} \otimes W$$

$$C_{o4} = D \otimes C_{w2}$$

$$C_{o5} = C_{d1} \otimes C_{w1}$$

$$C_{o6} = C_{d1} \otimes C_{w2}$$

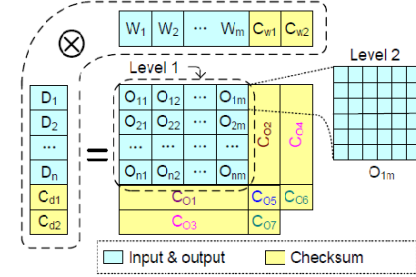
$$C_{o7} = C_{d2} \otimes C_{w1}$$

$$C_{d1} = \sum_{n=0}^{N-1} D_n$$

$$C_{d2} = \sum_{n=0}^{N-1} n D_n$$

$$C_{w1} = \sum_{m=0}^{M-1} W_m$$

$$C_{w2} = \sum_{m=0}^{M-1} m W_m$$



The row checksums used in this scheme are  $C_{o1}$  and  $C_{o3}$ .  $C_{o3}$  is computed from convolution operation between  $C_{d2}$  and  $W$ , and the related output summation is defined by  $S_{o3}[m] = \sum_{n=0}^{N-1} n \times O_{nm}$ .

For the detection of soft errors, we need to compare  $C_{o1}$  with  $S_{o1}$ . If they are not equal to each other at location  $j$ , the error can be located by  $i = \frac{C_{o3}[j] - S_{o3}[j]}{C_{o1}[j] - S_{o1}[j]}$  and  $j$ , and it can be corrected by adding  $C_{o1}[j] - S_{o1}[j]$  to the block  $(i, j)$ .

# Column Checksum Scheme (CIC)

$$C_{o1} = C_{d1} \otimes W$$

$$C_{o2} = D \otimes C_{w1}$$

$$C_{o3} = C_{d2} \otimes W$$

$$C_{o4} = D \otimes C_{w2}$$

$$C_{o5} = C_{d1} \otimes C_{w1}$$

$$C_{o6} = C_{d1} \otimes C_{w2}$$

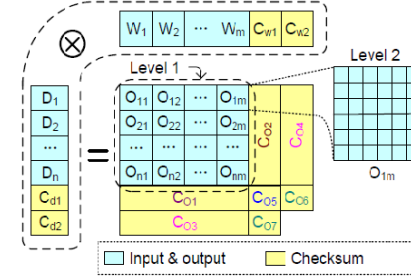
$$C_{o7} = C_{d2} \otimes C_{w1}$$

$$C_{d1} = \sum_{n=0}^{N-1} D_n$$

$$C_{d2} = \sum_{n=0}^{N-1} n D_n$$

$$C_{w1} = \sum_{m=0}^{M-1} W_m$$

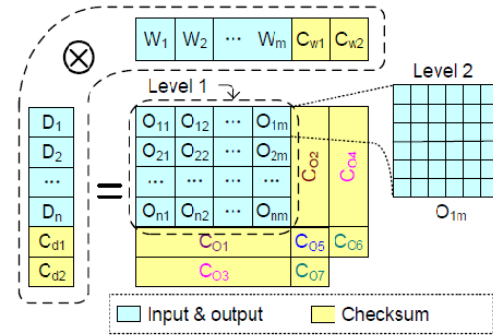
$$C_{w2} = \sum_{m=0}^{M-1} m W_m$$



The column checksums used in this scheme are  $C_{o2}$  and  $C_{o4}$ .  $C_{o4}$  is defined by performing convolution operation between  $D$  and  $C_{w2}$ , and the related output summation is defined as  $S_{o4}[n] = \sum_{m=0}^{M-1} m \times O_{nm}$ . To detect soft errors, we compare  $C_{o2}$  with  $S_{o2}$  first. If they are not equal to each other at location  $i$ , the error can be located by  $i$  and  $j$  ( $= \frac{C_{o4}[i] - S_{o4}[i]}{C_{o2}[i] - S_{o2}[i]}$ ), and it can be recovered by adding  $C_{o2}[i] - S_{o2}[i]$  to the block  $(i, j)$ .

# Checksum-of-Checksum Scheme (CoC/CoC-D)

$$\begin{aligned}
 C_{o1} &= C_{d1} \otimes W \\
 C_{o2} &= D \otimes C_{w1} \\
 C_{o3} &= C_{d2} \otimes W \\
 C_{o4} &= D \otimes C_{w2} \\
 C_{o5} &= C_{d1} \otimes C_{w1} \\
 C_{o6} &= C_{d1} \otimes C_{w2} \\
 C_{o7} &= C_{d2} \otimes C_{w1}
 \end{aligned}$$



$$\begin{aligned}
 C_{o5} &= \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} O_{nm} = S_{o5} \\
 C_{o6} &= \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} m \times O_{nm} = S_{o6} \\
 C_{o7} &= \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} n \times O_{nm} = S_{o7}
 \end{aligned}$$

- scheme involves neither D nor W but only their checksums, so it is named checksum-of-checksum scheme (or CoC scheme for short).

## Analysis of Protection Ability for Convolution Checksum Schemes

# Fault Model

## Assumptions:

- transient faults in computational units
  - data corruption faults (both transient and persistent) in memory (including cache).
- 
- Fault represents a malfunction event, and we denote its corresponding symptom as soft error.
  - **Soft error protection includes error detection and error correction.**
  - Error detection means that the scheme can detect soft errors without knowing the exact location. Error correction means that the scheme can locate the soft error locations and recover the incorrect result.

# Limitation

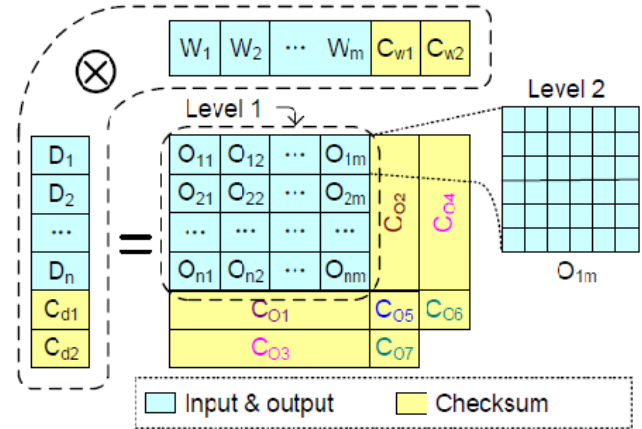
- **Without loss of generality, the following analysis considers at most one fault per convolution.**
- One convolutional neural network contains several or even tens of convolutional layers, and the total forward execution time of a CNN model is usually within seconds.
- Thus, we can reasonably assume that at most one fault may strike to one convolutional layer, considering the short executing time of a single layer.
- Multiple faults per convolution can also be detected by our schemes and recovered by recomputing the corrupted convolutional layer.

# Analysis of Soft Error in D and W

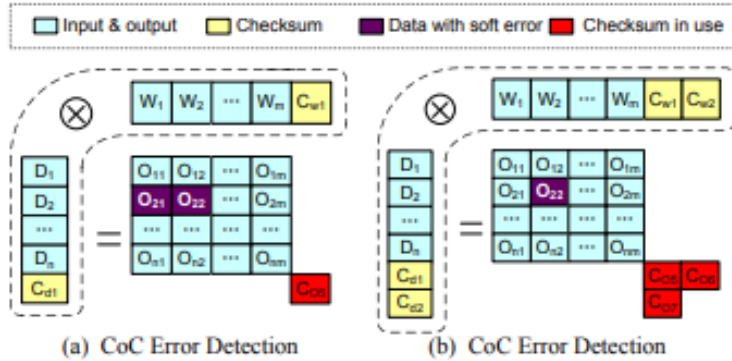
- One fault occurring during the convolution execution can result in multiple soft errors in W and D.
- **The soft errors in W can be detected by comparing the checksum of W with Cw1 and corrected by reloading weights from the CNN model.**
- **The soft errors in D do not need correction because D will be discarded after convolution computation.**
- The resulting errors in the output can be detected and corrected by the checksums of the output, as demonstrated below.



- **One fault during the convolution execution can result in corruption of one block row or column of O.**
- By definition, the row  $i$  of  $O$  is computed by the  $i$ th block of  $D$  with  $W$ . Thus, one fault in  $D$  would result in at most one corrupted row.
- The column  $j$  of  $O$  is computed by  $D$  with the  $j$ th block of  $W$ . Thus, one fault in  $W$  would result in at most one corrupted column.
- Moreover, the intermediate result will be reused only by the same row or column, such that one fault in the computational units would corrupt only values in the same row or column.
- Accordingly, the soft error protection ability in the context of at most one corrupted row or column of  $O$  will be decided.



# Soft Error Protection Ability of CoC Scheme

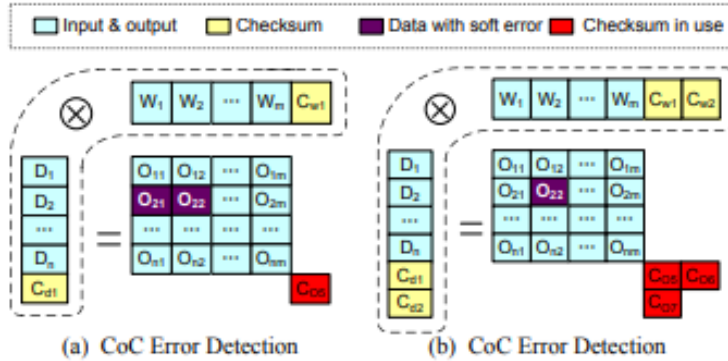


$$\begin{aligned}
 C_{o1} &= C_{d1} \otimes W \\
 C_{o2} &= D \otimes C_{w1} \\
 C_{o3} &= C_{d2} \otimes W \\
 C_{o4} &= D \otimes C_{w2} \\
 C_{o5} &= C_{d1} \otimes C_{w1} \\
 C_{o6} &= C_{d1} \otimes C_{w2} \\
 C_{o7} &= C_{d2} \otimes C_{w1}
 \end{aligned}$$

When soft errors strike the input or output data:

**Multiple soft errors can be detected by using only Co5.**

# Soft Error Protection Ability of CoC Scheme



$$\begin{aligned}
 C_{o1} &= C_{d1} \otimes W \\
 C_{o2} &= D \otimes C_{w1} \\
 C_{o3} &= C_{d2} \otimes W \\
 C_{o4} &= D \otimes C_{w2} \\
 C_{o5} &= C_{d1} \otimes C_{w1} \\
 C_{o6} &= C_{d1} \otimes C_{w2} \\
 C_{o7} &= C_{d2} \otimes C_{w1}
 \end{aligned}$$

When soft errors strike the input or output data:

**A single soft error in O can be corrected by CoC using all checksums including Co5, Co6, and Co7, as shown in Figure 2 (b).**

# Soft Error Protection Ability of CoC Scheme

$$C_{o5} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} O_{nm} = S_{o5}$$

$$C_{o6} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} m \times O_{nm} = S_{o6}$$

$$C_{o7} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} n \times O_{nm} = S_{o7}$$

$$C_{o5} - S_{o5} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} O'_{nm} - O_{nm} = \delta$$

$$C_{o6} - S_{o6} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} m \times (O'_{nm} - O_{nm}) = j \times \delta$$

$$C_{o7} - S_{o7} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} n \times (O'_{nm} - O_{nm}) = i \times \delta$$

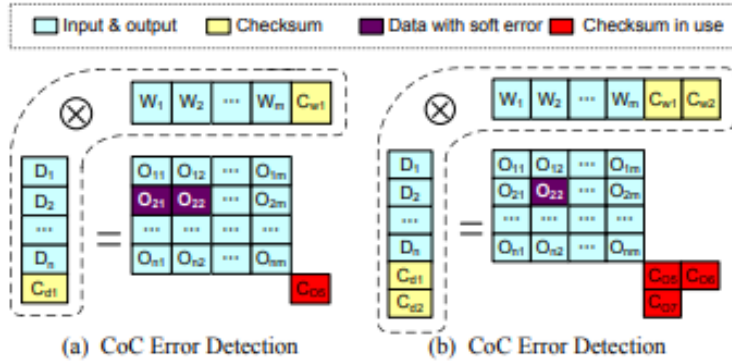
The location  $i, j$  can be obtained by  $i = (C_{o7} - S_{o7})/\delta$  and  $j = (C_{o6} - S_{o6})/\delta$ . Then the soft error can be fixed by adding  $\delta$  to  $O_{ij}$ .

When soft errors strike the input or output data:

**A single soft error in  $O$  can be corrected by CoC using all checksums including  $Co5$ ,  $Co6$ , and  $Co7$ .**

**[CoC cannot correct soft errors across multiple blocks in  $O$ .]**

# Soft Error Protection Ability of CoC Scheme



$$\begin{aligned}
 C_{o1} &= C_{d1} \otimes W \\
 C_{o2} &= D \otimes C_{w1} \\
 C_{o3} &= C_{d2} \otimes W \\
 C_{o4} &= D \otimes C_{w2} \\
 C_{o5} &= C_{d1} \otimes C_{w1} \\
 C_{o6} &= C_{d1} \otimes C_{w2} \\
 C_{o7} &= C_{d2} \otimes C_{w1}
 \end{aligned}$$

When soft errors strike the input or output data:

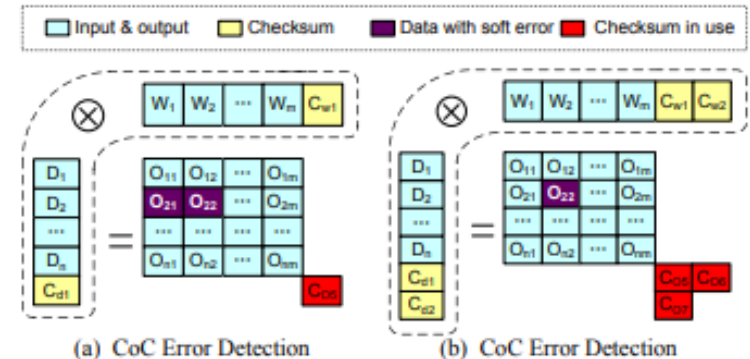
**However, CoC cannot correct soft errors across multiple blocks in O.**

When checksums are corrupted

Such soft errors can cause inconsistency among the output checksums of CoC, which can be used for error detection.

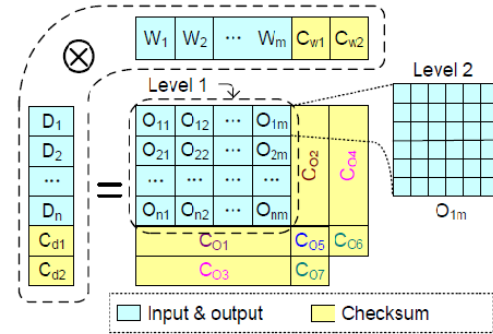
- For example, if  $C_{d1}$  is corrupted, leading to corrupted  $C_{o5}$  and  $C_{o6}$  with correct  $C_{o7}$ .
- **We can detect this abnormal pattern when comparing checksums with the summation of O to detect the input checksum corruption.**
- The input D, W, and output O are clean and without soft errors since fault frequency is at most once per convolution. Thus, we can safely discard all the checksums and finish this convolution computation

$$\begin{aligned} C_{o1} &= C_{d1} \otimes W \\ C_{o2} &= D \otimes C_{w1} \\ C_{o3} &= C_{d2} \otimes W \\ C_{o4} &= D \otimes C_{w2} \\ C_{o5} &= C_{d1} \otimes C_{w1} \\ C_{o6} &= C_{d1} \otimes C_{w2} \\ C_{o7} &= C_{d2} \otimes C_{w1} \end{aligned}$$



## When checksums are corrupted

$$\begin{aligned}
 C_{o1} &= C_{d1} \otimes W \\
 C_{o2} &= D \otimes C_{w1} \\
 C_{o3} &= C_{d2} \otimes W \\
 C_{o4} &= D \otimes C_{w2} \\
 C_{o5} &= C_{d1} \otimes C_{w1} \\
 C_{o6} &= C_{d1} \otimes C_{w2} \\
 C_{o7} &= C_{d2} \otimes C_{w1}
 \end{aligned}$$



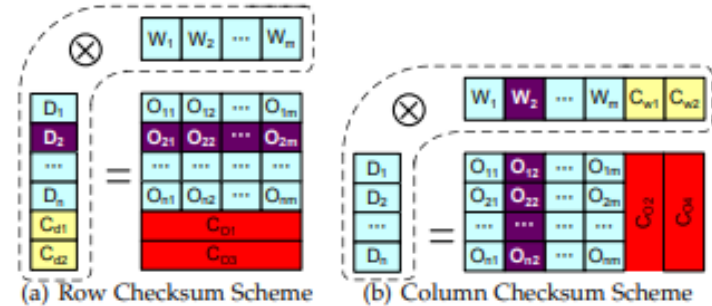
$$\begin{aligned}
 C_{o5} &= \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} O_{nm} = S_{o5} \\
 C_{o6} &= \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} m \times O_{nm} = S_{o6} \\
 C_{o7} &= \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} n \times O_{nm} = S_{o7}
 \end{aligned}$$

- Only one of them will be corrupted among **Co5**, **Co6**, and **Co7**

[ 'One fault during the convolution execution': Assumption from fault model]

# Ability of Row and Column Checksum Scheme

- the row checksum scheme and column checksum scheme are symmetric.
- the row checksum scheme can detect and correct soft errors if they are in the same row.**
- If the soft errors are in the same column, the row checksum scheme can only detect soft errors; it has no correction ability.**
- The column checksum scheme, on the contrary, can detect and correct errors located in the same column but fail to correct those appearing in the same row.





# Revisit Row Checksum Scheme (RC)

$$C_{o1} = C_{d1} \otimes W$$

$$C_{o2} = D \otimes C_{w1}$$

$$C_{o3} = C_{d2} \otimes W$$

$$C_{o4} = D \otimes C_{w2}$$

$$C_{o5} = C_{d1} \otimes C_{w1}$$

$$C_{o6} = C_{d1} \otimes C_{w2}$$

$$C_{o7} = C_{d2} \otimes C_{w1}$$

The row checksums used in this scheme are  $C_{o1}$  and  $C_{o3}$ .  $C_{o3}$  is computed from convolution operation between  $C_{d2}$  and  $W$ , and the related output summation is defined by  $S_{o3}[m] = \sum_{n=0}^{N-1} n \times O_{nm}$ .

For the detection of soft errors, we need to compare  $C_{o1}$  with  $S_{o1}$ . If they are not equal to each other at location  $j$ , the error can be located by  $i = \frac{C_{o3}[j] - S_{o3}[j]}{C_{o1}[j] - S_{o1}[j]}$  and  $j$ , and it can be corrected by adding  $C_{o1}[j] - S_{o1}[j]$  to the block  $(i, j)$ .

# Revisit Column Checksum Scheme (CIC)

$$C_{o1} = C_{d1} \otimes W$$

$$C_{o2} = D \otimes C_{w1}$$

$$C_{o3} = C_{d2} \otimes W$$

$$C_{o4} = D \otimes C_{w2}$$

$$C_{o5} = C_{d1} \otimes C_{w1}$$

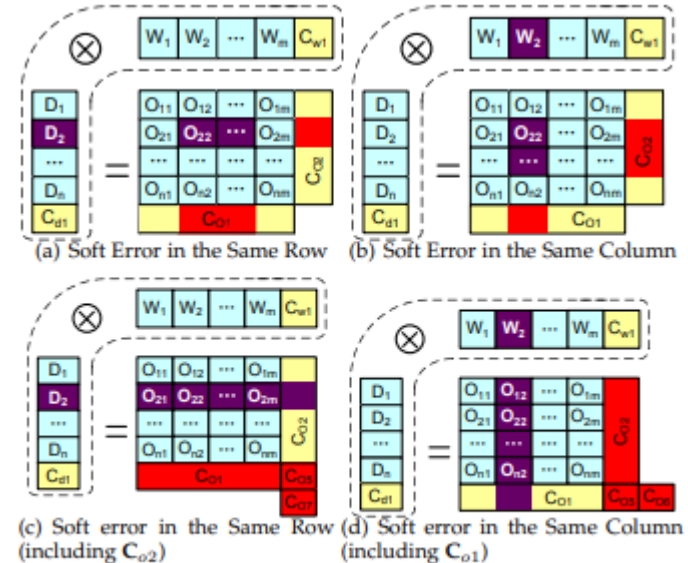
$$C_{o6} = C_{d1} \otimes C_{w2}$$

$$C_{o7} = C_{d2} \otimes C_{w1}$$

The column checksums used in this scheme are  $C_{o2}$  and  $C_{o4}$ .  $C_{o4}$  is defined by performing convolution operation between  $D$  and  $C_{w2}$ , and the related output summation is defined as  $S_{o4}[n] = \sum_{m=0}^{M-1} m \times O_{nm}$ . To detect soft errors, we compare  $C_{o2}$  with  $S_{o2}$  first. If they are not equal to each other at location  $i$ , the error can be located by  $i$  and  $j$  ( $= \frac{C_{o4}[i] - S_{o4}[i]}{C_{o2}[i] - S_{o2}[i]}$ ), and it can be recovered by adding  $C_{o2}[i] - S_{o2}[i]$  to the block  $(i, j)$ .

# Ability of Full Checksum Scheme

- Full checksum scheme has the highest ability to correct soft errors.
- The scheme uses both the row checksum  $Co1$  and column checksum  $Co2$  so that it can correct soft errors in both directions.



Soft Error Protection Ability of Full Checksum Scheme

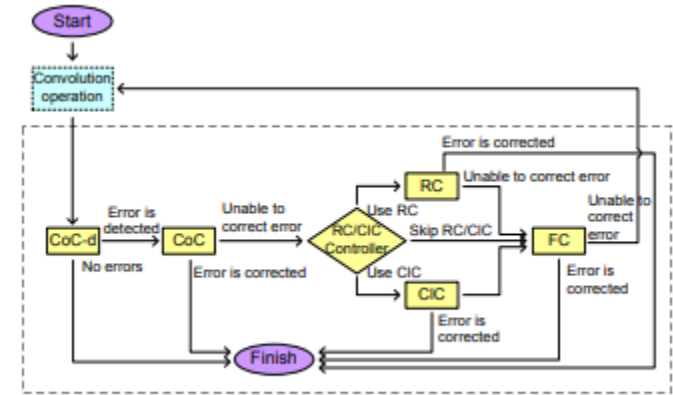
## Ability of four schemes

- CoC scheme has the lowest error correction ability and that the full checksum scheme has the best error correction ability.
- The abilities of the row checksum scheme and column checksum scheme are higher than that of the CoC scheme but lower than that of the full checksum scheme.
- CoC-D can detect multiple soft errors but without correction ability.

**The analysis here serves as the fundamental basis of our low-overhead high protection design.**

# Multischeme Workflow for Soft Error Protection

- To achieve the highest protection ability and lowest overhead: multischeme workflow by integrating the four schemes.
- use CoC-D to detect errors because it has the lowest overhead.
- For the error correction, put CoC in the beginning because it is the most lightweight method.
- By comparison, FC has highest correction ability but also highest time overhead, so put it at the end of the workflow.



Multischeme Workflow Designed to Detect/Correct Soft Errors

# Soft Error Protection Ability of CoC Scheme

$$C_{o5} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} O_{nm} = S_{o5}$$

$$C_{o6} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} m \times O_{nm} = S_{o6}$$

$$C_{o7} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} n \times O_{nm} = S_{o7}$$

$$C_{o5} - S_{o5} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} O'_{nm} - O_{nm} = \delta$$

$$C_{o6} - S_{o6} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} m \times (O'_{nm} - O_{nm}) = j \times \delta$$

$$C_{o7} - S_{o7} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} n \times (O'_{nm} - O_{nm}) = i \times \delta$$

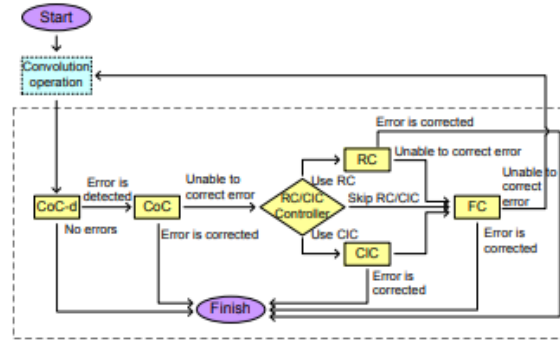
The location  $i, j$  can be obtained by  $i = (C_{o7} - S_{o7})/\delta$  and  $j = (C_{o6} - S_{o6})/\delta$ . Then the soft error can be fixed by adding  $\delta$  to  $O_{ij}$ .

When soft errors strike the input or output data:

**A single soft error in  $O$  can be corrected by CoC using all checksums including  $Co5$ ,  $Co6$ , and  $Co7$ .**

**[CoC cannot correct soft errors across multiple blocks in  $O$ .]**

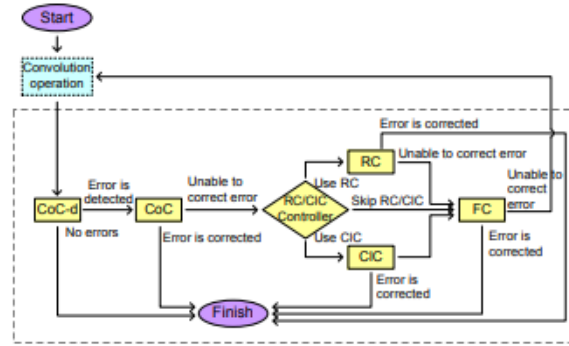
# Multischeme Workflow for Soft Error Protection



Multischeme Workflow Designed to Detect/Correct Soft Errors

- **The error detection modules will be executed for every execution whether there is a soft error or not.** Thus, any unnecessary computations should be avoided in order to reduce the overall overhead.
- For instance, both CoC-D and FC are able to detect all the soft errors, but we adopt only CoC-D in the workflow for error detection because FC has a much higher overhead. RC and CIC cannot detect soft errors correctly if the checksum is corrupted.

# Multischeme Workflow for Soft Error Protection

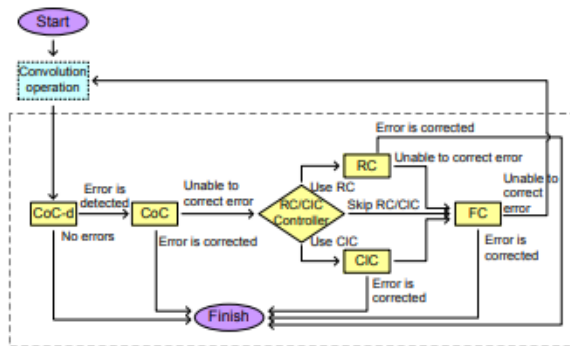


Multischeme Workflow Designed to Detect/Correct Soft Errors

- **The error correction module will not be executed until some soft errors are detected.** The schemes in this module will be invoked to fix soft errors according to the workflow. If it fails to correct the errors due to inconsistency of checksum blocks or illegal error locations, the next-level scheme will be invoked.



# Multischeme Workflow for Soft Error Protection



Multischeme Workflow Designed to Detect/Correct Soft Errors

- **Checksums can be reused among different CNN schemes in the workflow**, the runtime of the workflow is actually lower than the sum of all schemes' runtimes.
- For example, both CoC-D and CoC use Co5; if CoC-D detects soft errors and CoC is invoked to correct soft errors, CoC can save the time of computing Co5 and its corresponding summation So5, since they have been computed by CoC-D.

# Code Ref

<https://oneapi-src.github.io/oneDNN/v0/index.html>

[https://github.com/oneapi-src/oneDNN/blob/master/examples/cnn\\_inference\\_f32.cpp](https://github.com/oneapi-src/oneDNN/blob/master/examples/cnn_inference_f32.cpp)

[https://github.com/oneapi-src/oneDNN/blob/master/examples/cnn\\_inference\\_int8.cpp](https://github.com/oneapi-src/oneDNN/blob/master/examples/cnn_inference_int8.cpp)

# Assignment

- **Implement** Multischeme Workflow for Soft Error Protection
- **If soft errors exist in Co1/Co2 (any of them), error detection and correction is still possible.**
  - **Justify or Refute**

Ref:

[https://www.researchgate.net/publication/348144782\\_FT-CNN\\_Algorithm-Based\\_Fault\\_Tolerance\\_for\\_Convolutional\\_Neural\\_Networks](https://www.researchgate.net/publication/348144782_FT-CNN_Algorithm-Based_Fault_Tolerance_for_Convolutional_Neural_Networks)