# Development of C programs for Convolutional Neural Network Accelerators

JiYoung An, Sujin Kang

Prof. Nikil Dutt , Kenshu Seto, Hamid Nejatollahi

Dept. of Computer Engineering in Kyung Hee University

Dept. of Computer Engineering in Hanyang University

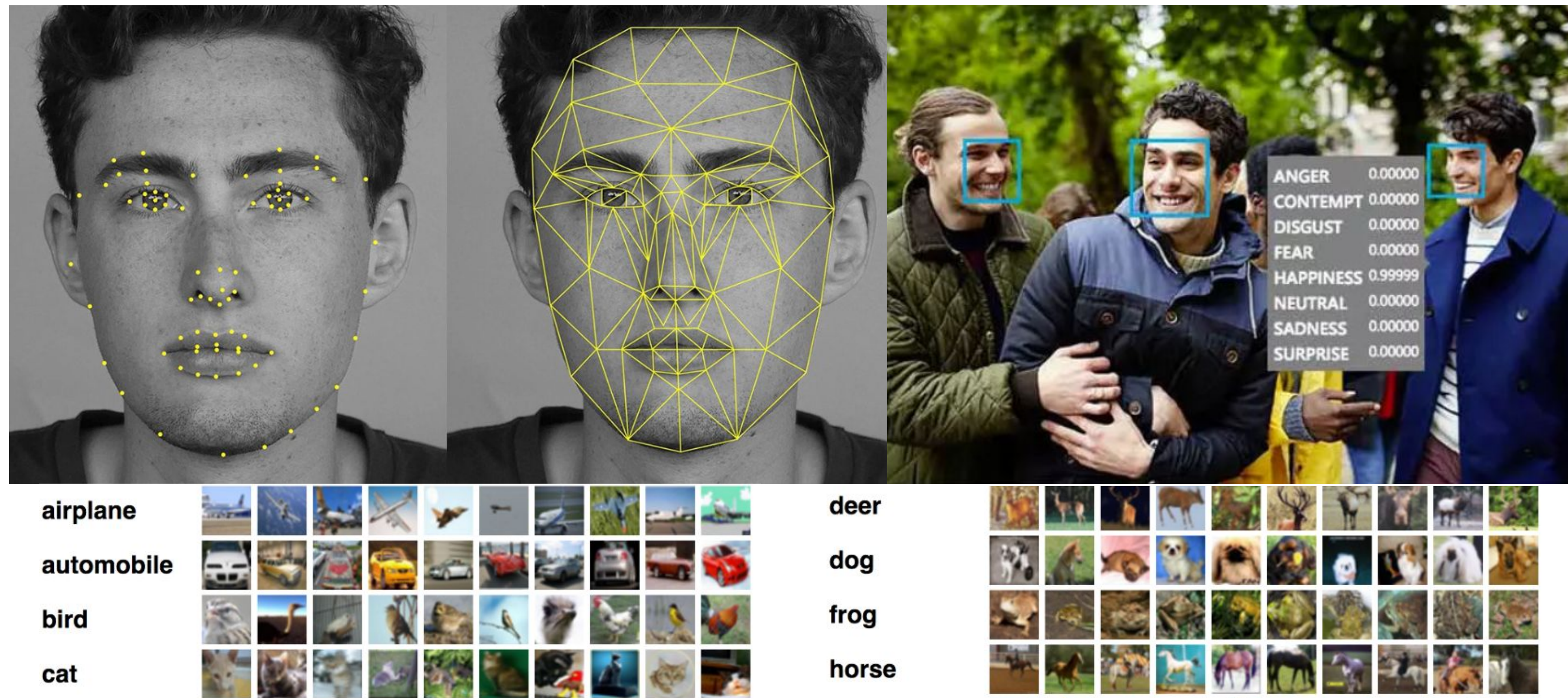University of California Irvine

University in Tokyo
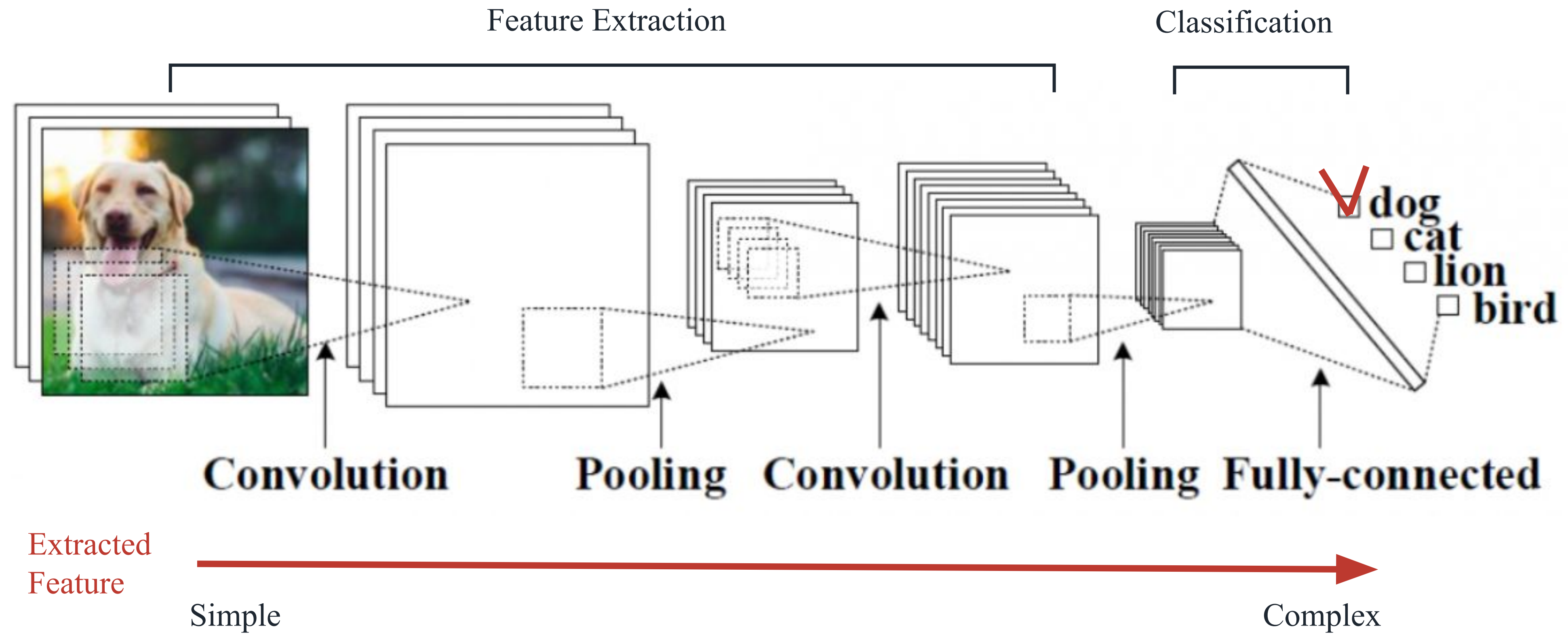
# BACKGROUND

# Convolution Neural Network

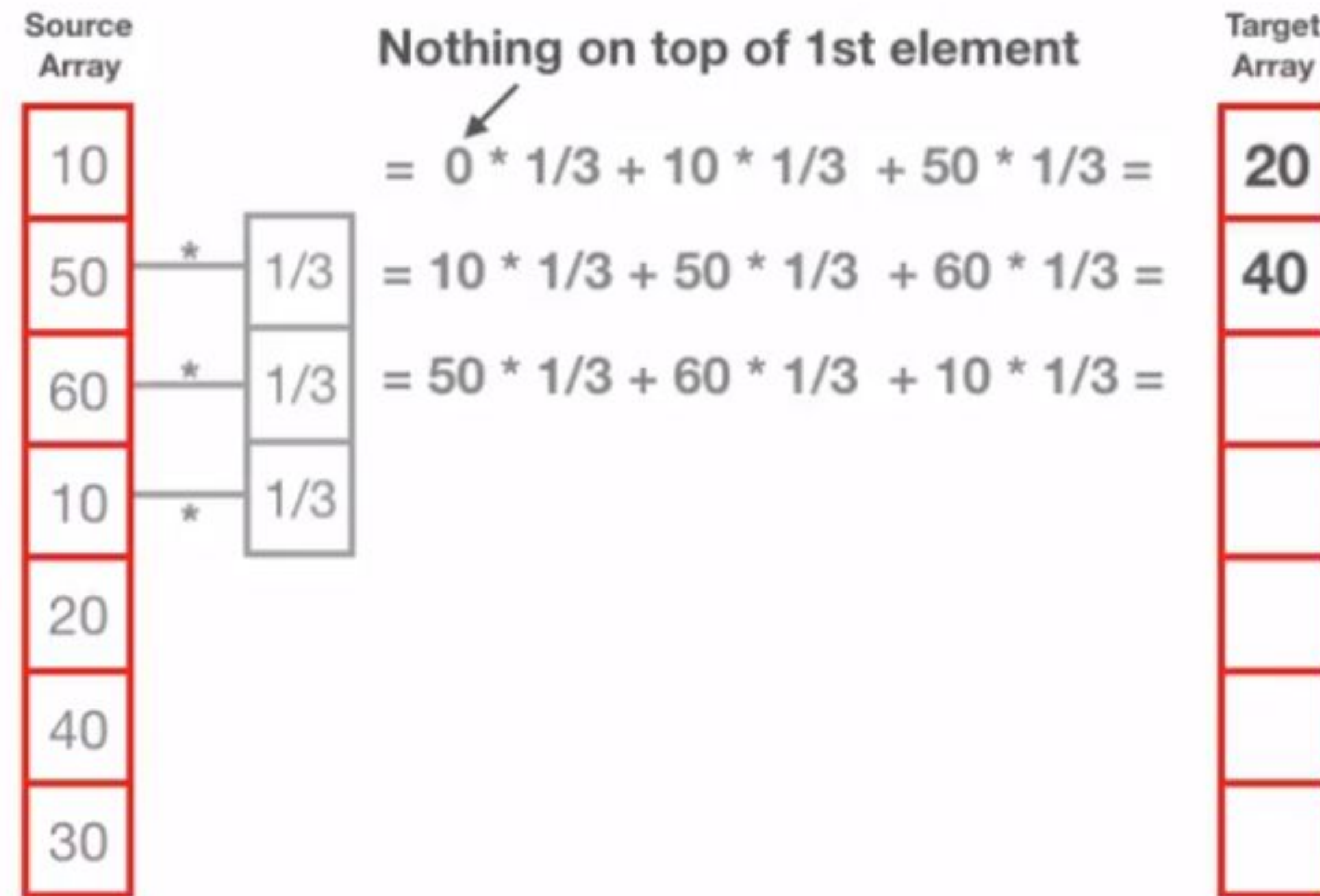Keyword : Massive data , Image recognition and classification
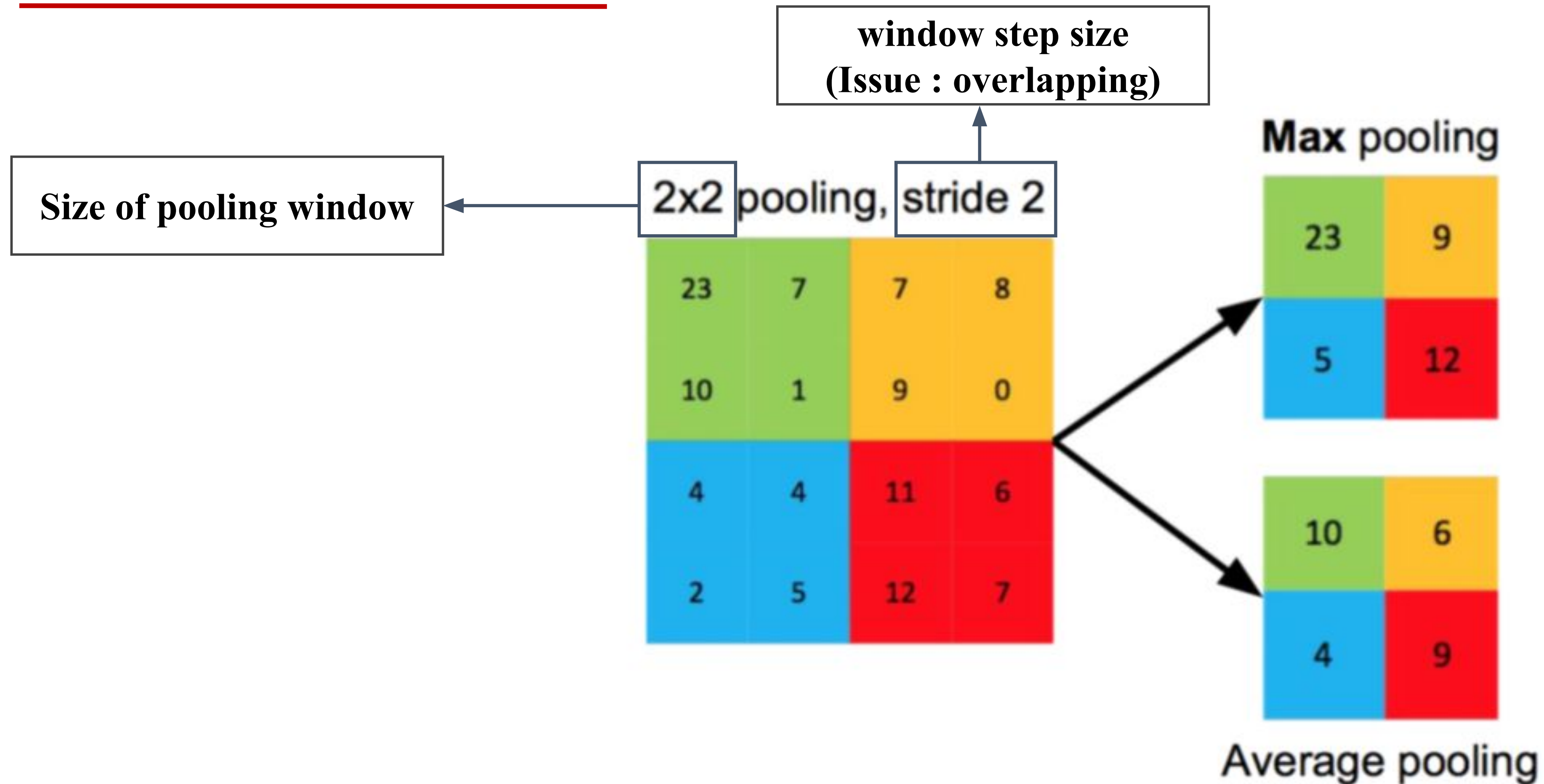
# 1. Convolutional Neural Networks



Keep spatial Information of Image

Automatically extract and learn its features

# 2. Convolution layer

Output = Source * Filter ( * Convolution )
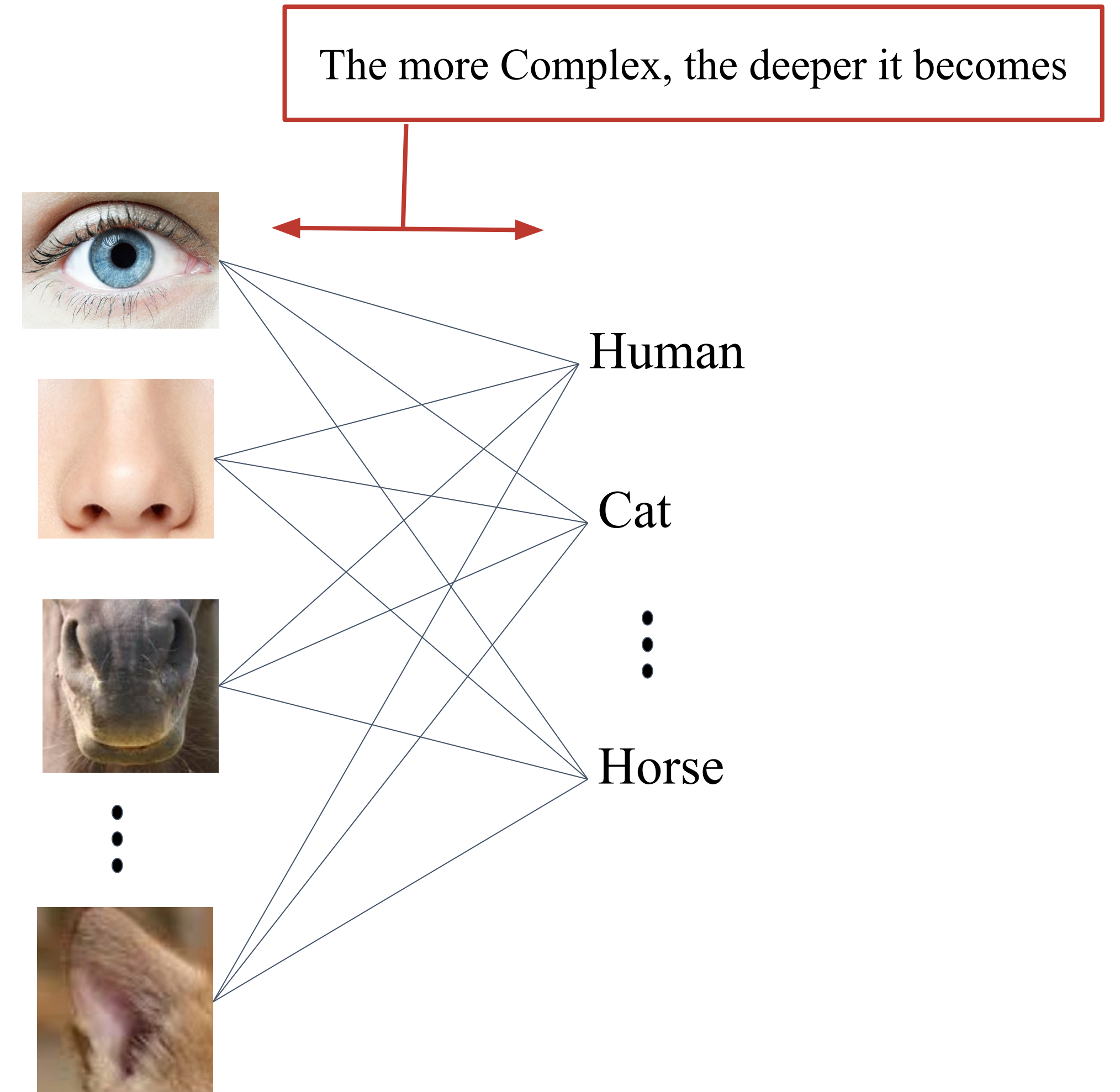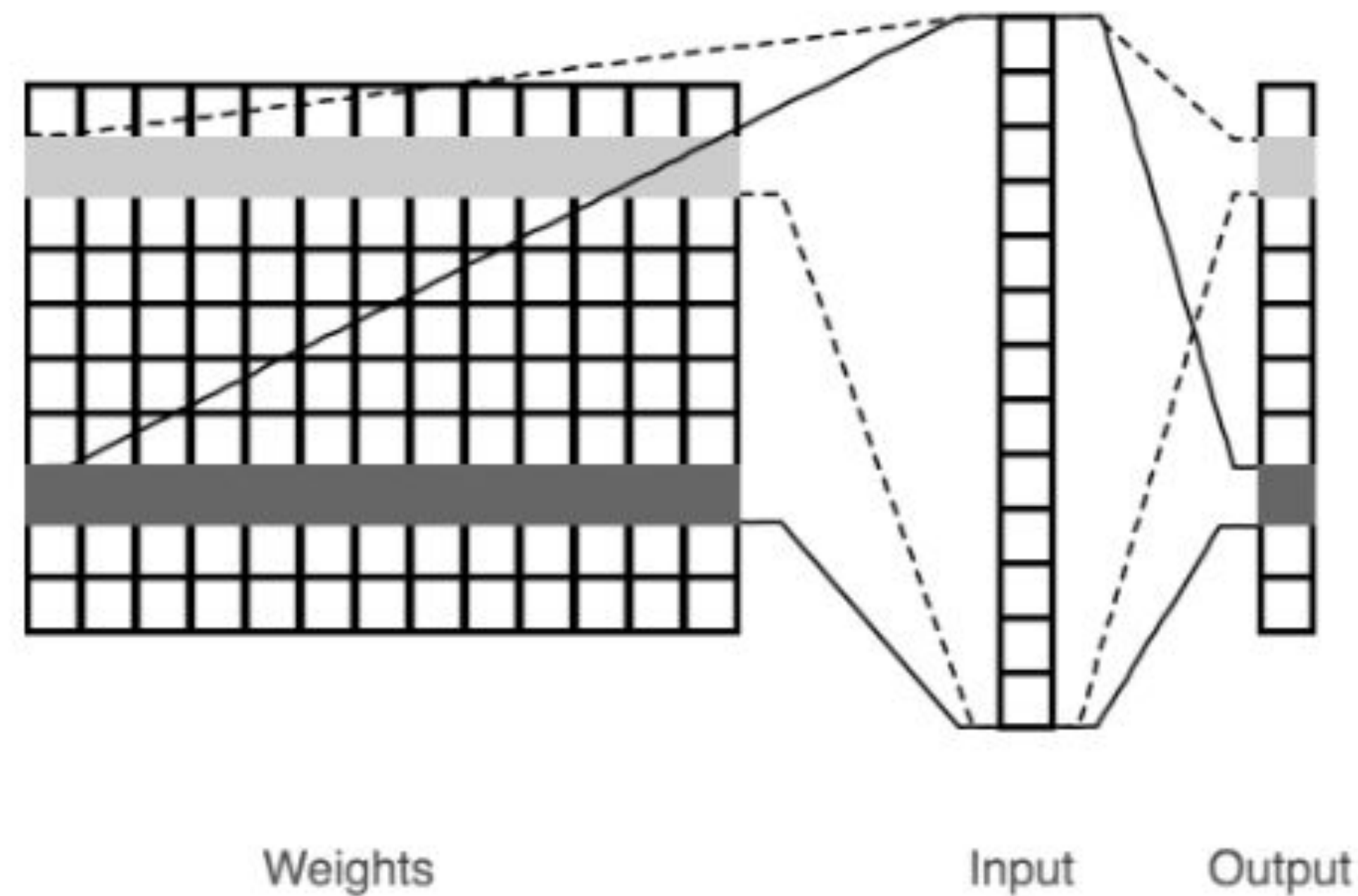
# 3. Pooling layer

**window step size (Issue : overlapping)**

**Size of pooling window**

2x2 pooling, stride 2

1. Extract strong feature and reduce resolution
2. Increase translation-invariance and noise-resilience

# 4. Fully connected layer



The more Complex, the deeper it becomes

Weights · Input · Output

Human
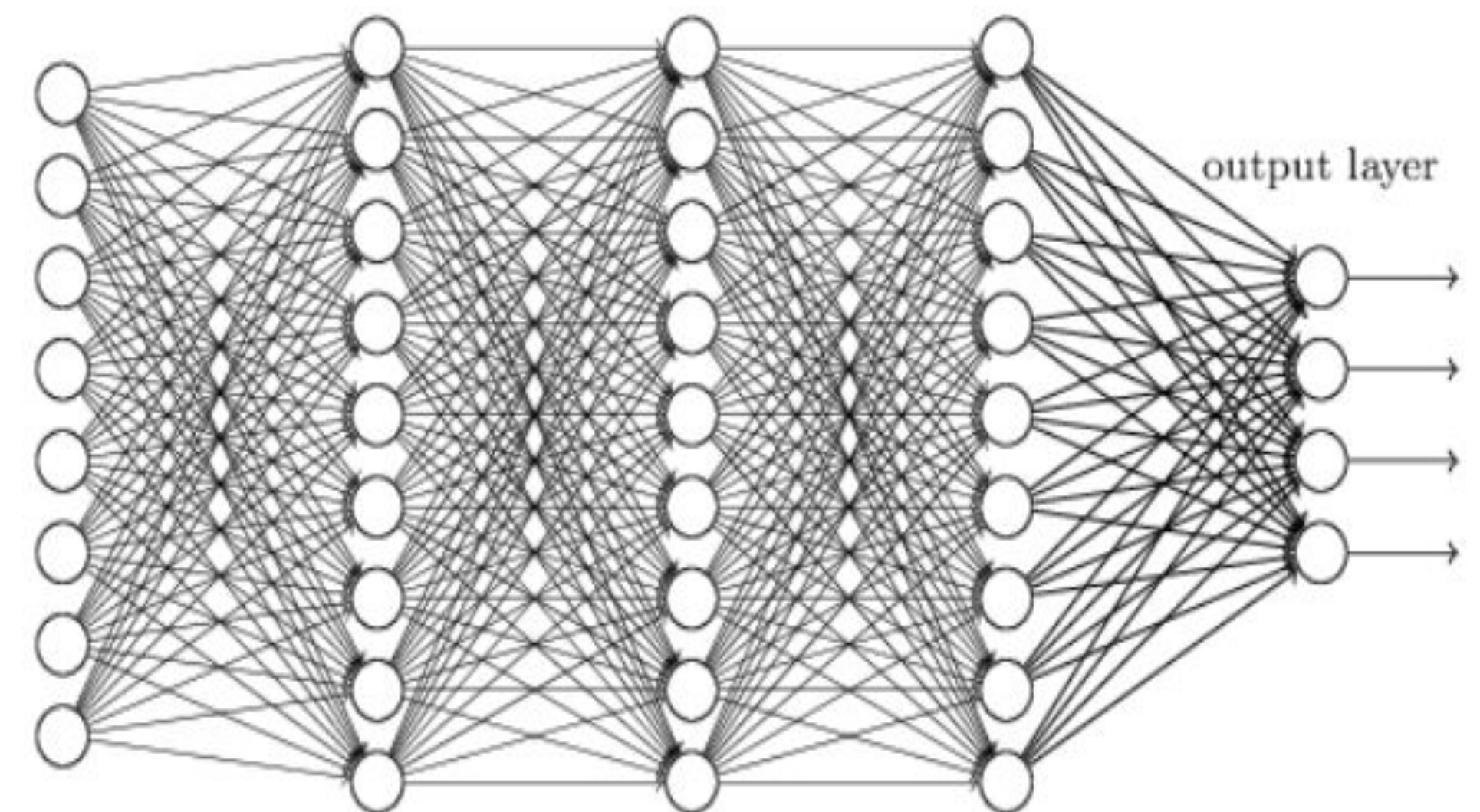
Cat

⋮

Horse

With Extracted Features, do classification

# PROJECT

# Problem of CNNs

Keyword : Time-consuming, Computing Energy, Calculation



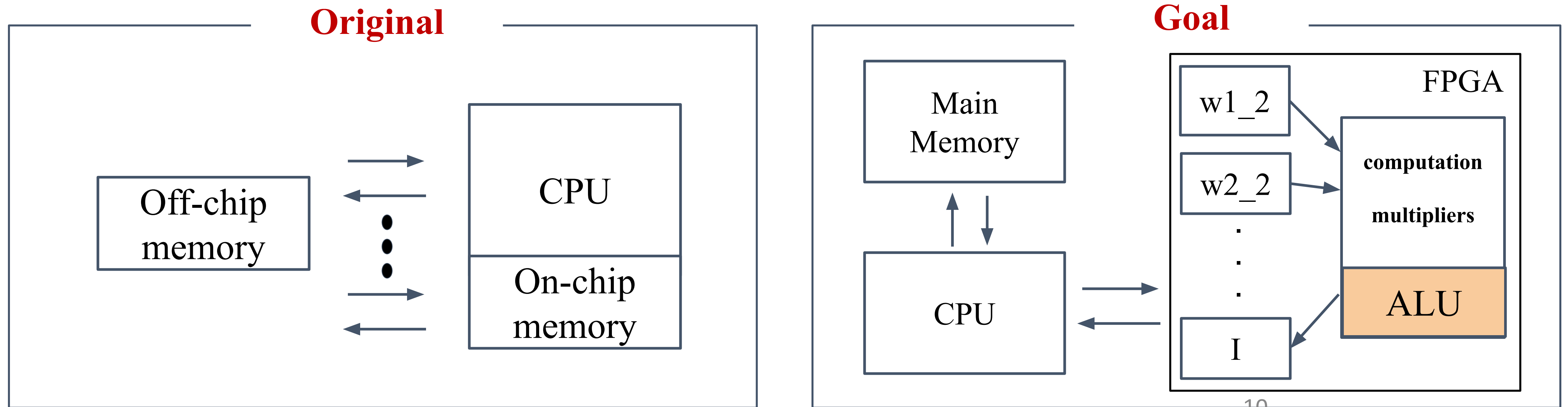| Layer (type) | Output Shape | Param # | Connected to |
|---|---|---|---|
| maxpooling2d_24 (MaxPooling2D) | (None, 512, 7, 7) | 0 | maxpooling2d_input_4[0][0] |
| batchnormalization_10 (BatchNorm | (None, 512, 7, 7) | 1024 | maxpooling2d_24[0][0] |
| flatten_8 (Flatten) | (None, 25088) | 0 | batchnormalization_10[0][0] |
| dense_22 (Dense) | (None, 4096) | 102764544 | flatten_8[0][0] |
| dropout_15 (Dropout) | (None, 4096) | 0 | dense_22[0][0] |
| batchnormalization_11 (BatchNorm | (None, 4096) | 8192 | dropout_15[0][0] |
| dense_23 (Dense) | (None, 4096) | 16781312 | batchnormalization_11[0][0] |
| dropout_16 (Dropout) | (None, 4096) | 0 | dense_23[0][0] |
| batchnormalization_12 (BatchNorm | (None, 4096) | 8192 | dropout_16[0][0] |
| dense_24 (Dense) | (None, 2) | 8194 | batchnormalization_12[0][0] |

Total params: 119571458
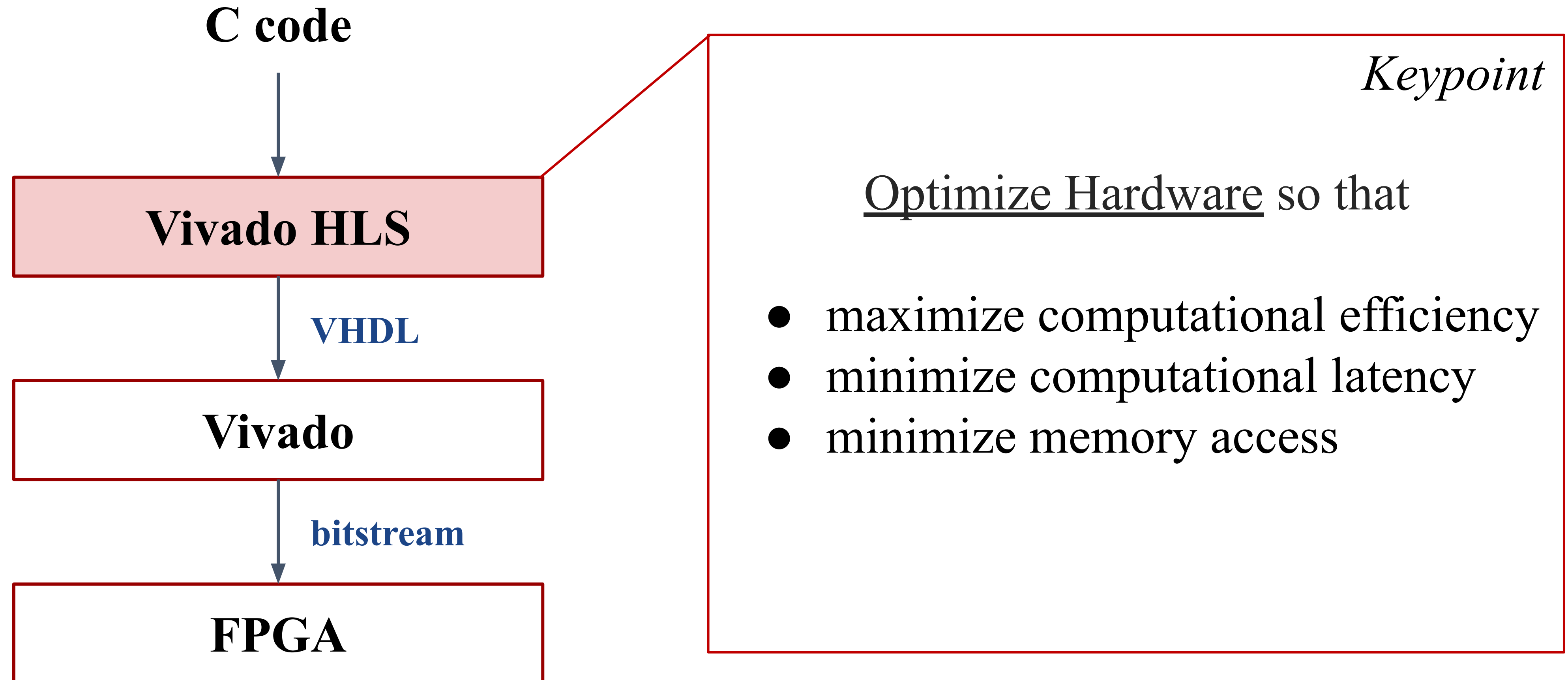
output layer

Required a lot of Computing resources and time

# Purpose

## Acceleration of Convolution Neural Networks to embed in light machine

Focus on Hardware Acceleration

# Approach

C code

Vivado HLS

VHDL

Vivado

bitstream

FPGA

*Keypoint*

Optimize Hardware so that

- maximize computational efficiency
- minimize computational latency
- minimize memory access

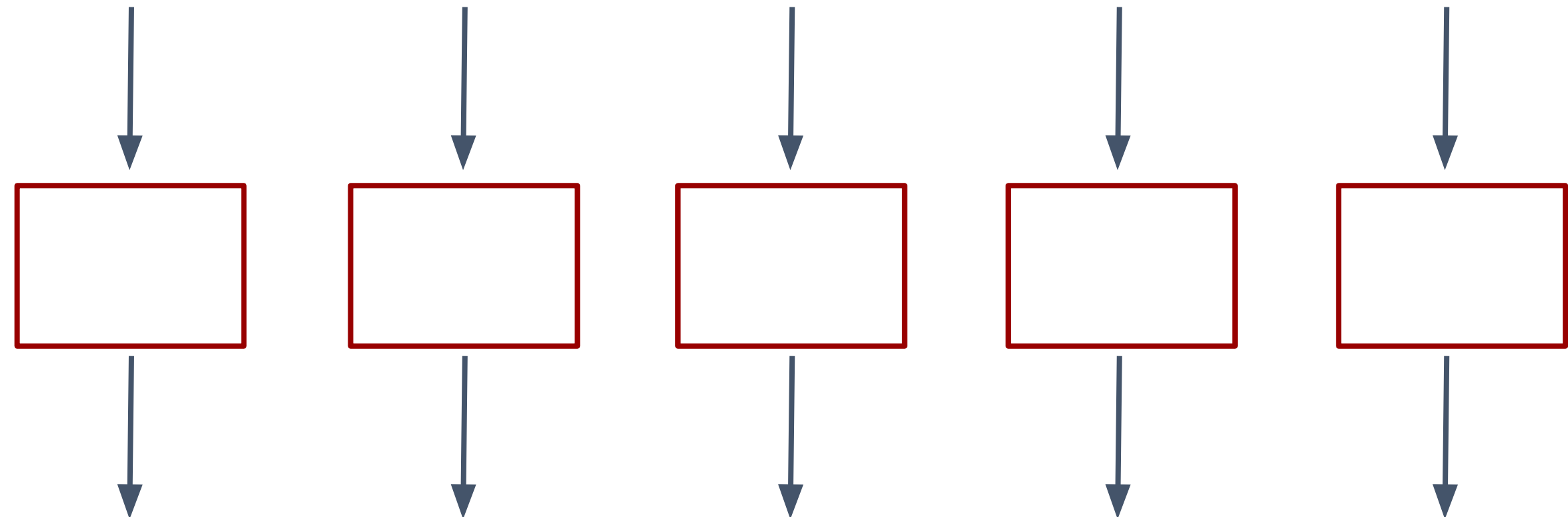# 1. Hardware Optimization : Array partition

**Normal**

**Array partition : Complete**

*Command Ex :*
*#pragma HLS ARRAY_PARTITION variable=B complete*

***Pros and Cons : Parallel access  vs Hardware delay***

# 2. Hardware Optimization : Pipeline

*In Convolution Code...*
for( m = 0 ; m< 4096 ; m ++)
  for( n = 0 ; n < 2048 ; n ++ )
    for( i = 0 ; i < 256 ; i ++ )
      for( j = 0 ; j < 256 ; j ++ ){
        for(x =0 ; x <3 ; x ++){
          for(y = 0 ; y <3 ; y ++ ){
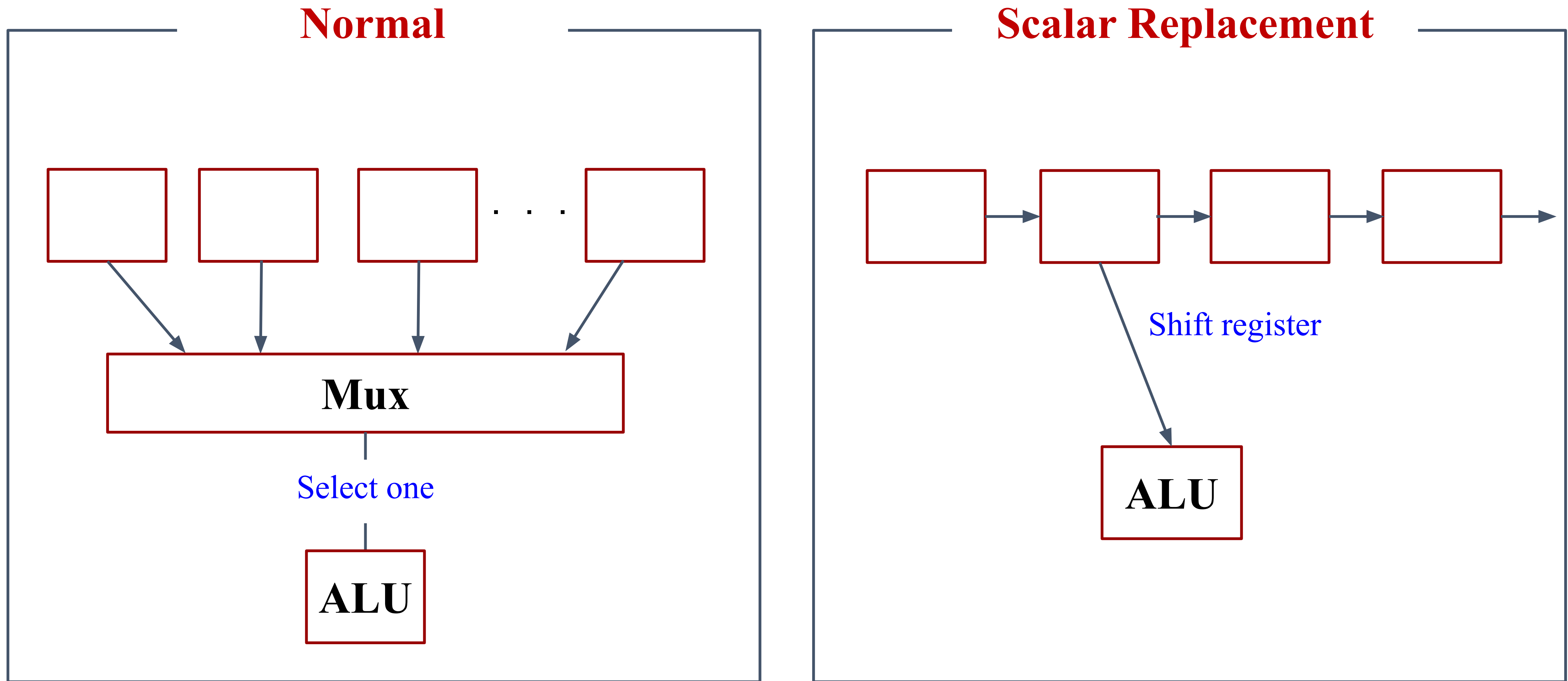            O[i][j]  += I[i][j] * W[m][n][x][y]
          }
        }
      }

**Apply Pipeline**

*Command Ex : #pragma HLS PIPELINE*

# Our Role

## 1. Template library

Implement function
in CNN model
without model dependency
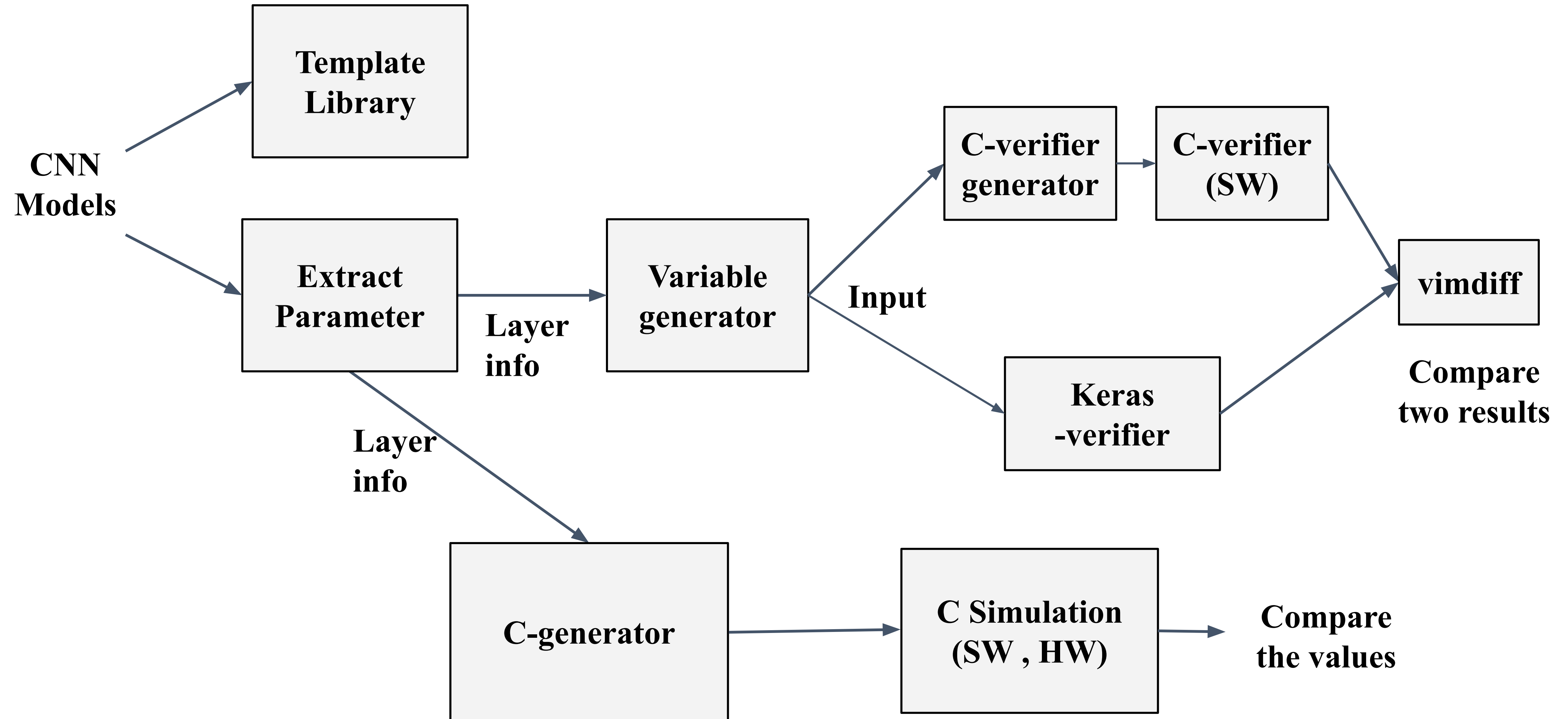
## 2. Verifier

SW version
vs
Python(keras)

## 3. Generator

Generate C code for Vivado
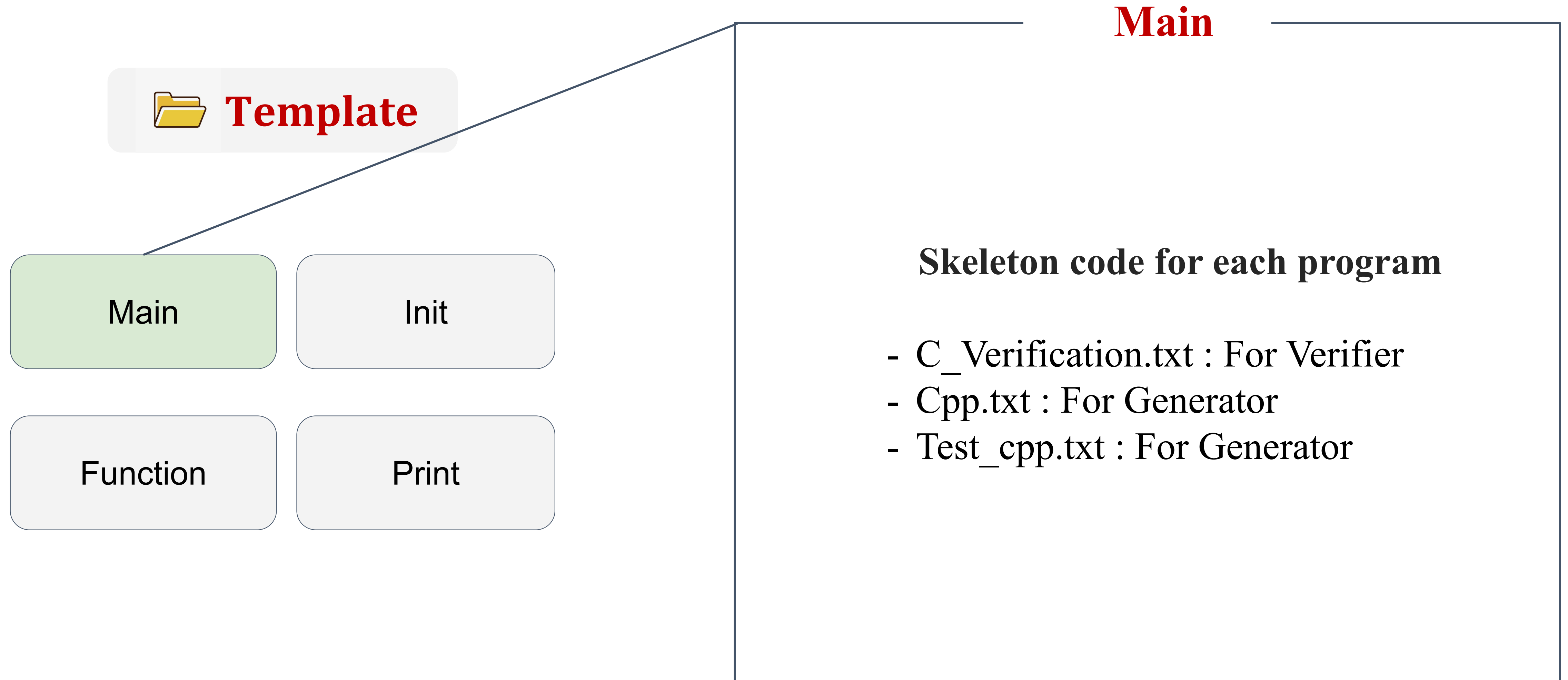
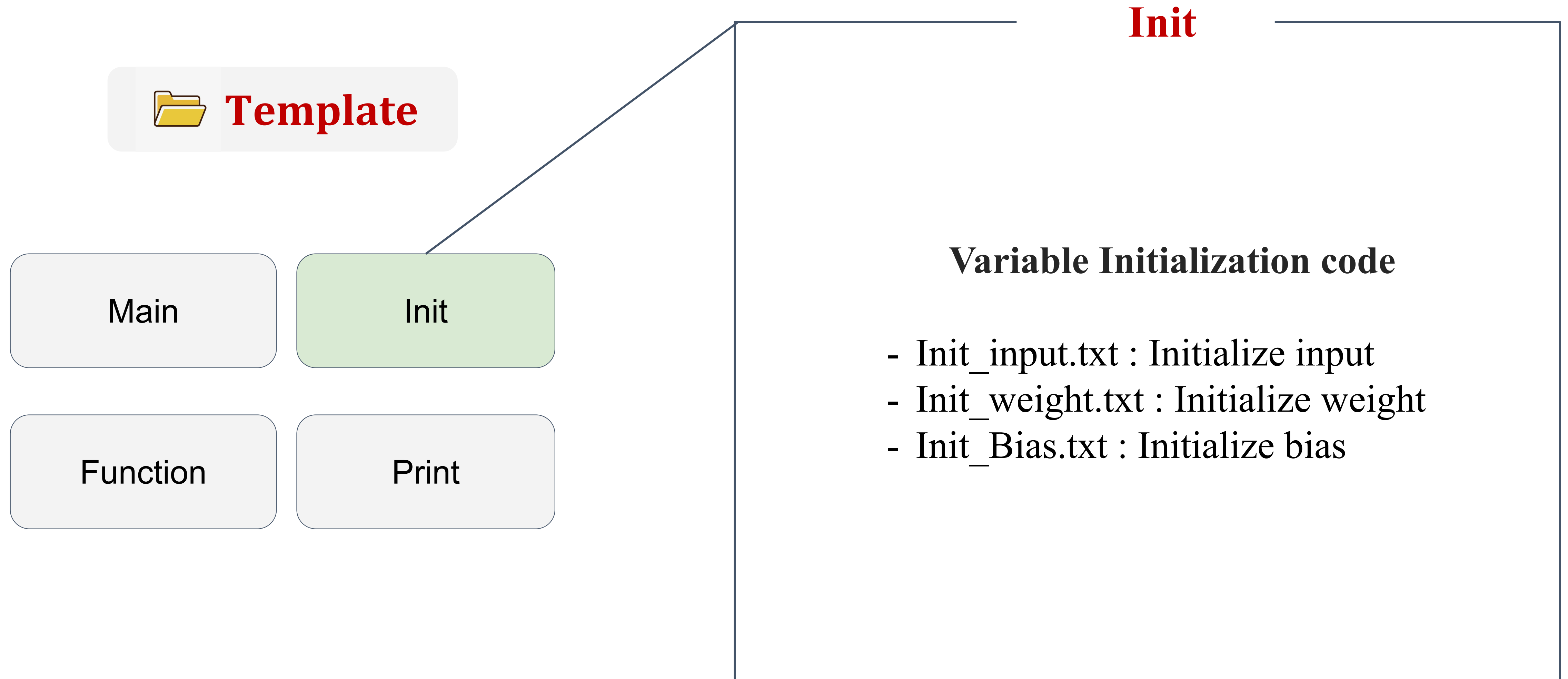## 4. C simulation

SW version
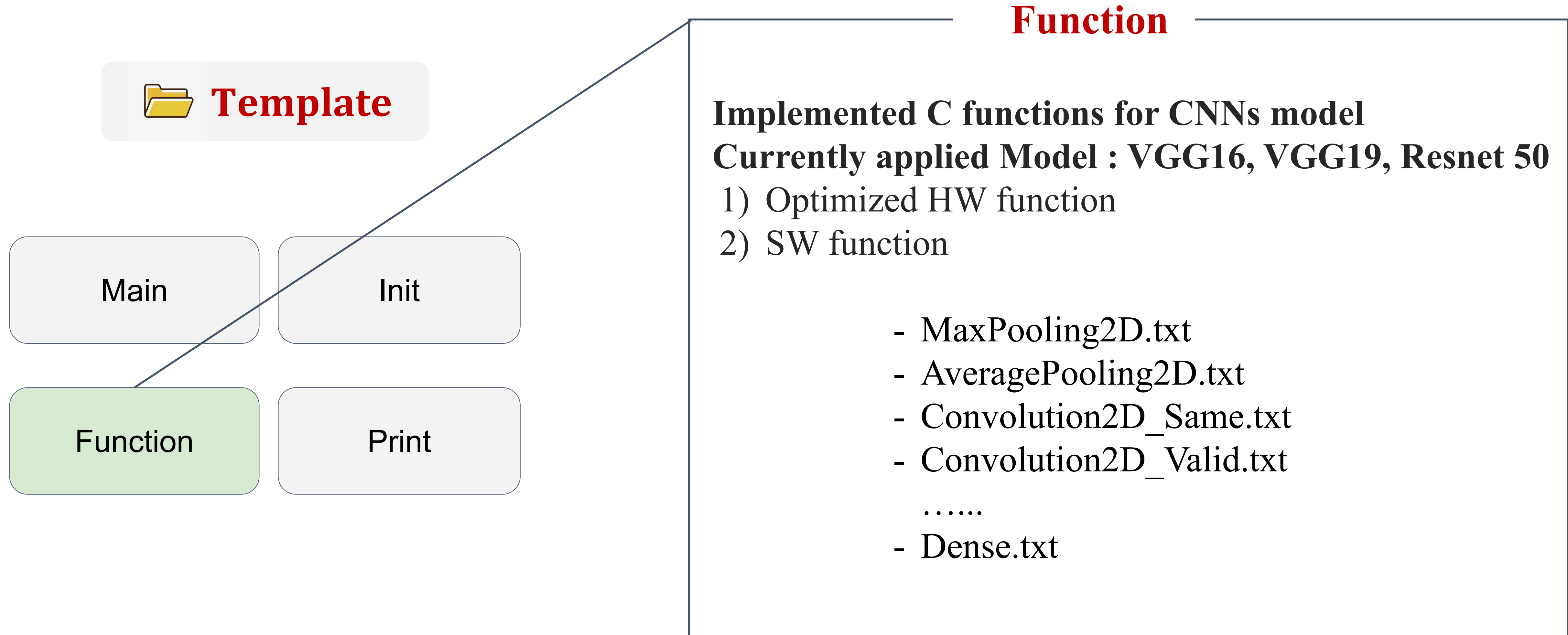vs
HW version

# Work Process

# 1. Template Library

📂 **Template**

| Main | Init |
|------|------|

| Function | Print |
|----------|-------|

**Main**

**Skeleton code for each program**

- C_Verification.txt : For Verifier
- Cpp.txt : For Generator
- Test_cpp.txt : For Generator

# 1. Template Library

📂 **Template**

| | |
|---|---|
| Main | **Init** |
| Function | Print |

## Init

### Variable Initialization code

- Init_input.txt : Initialize input
- Init_weight.txt : Initialize weight
- Init_Bias.txt : Initialize bias

# 1. Template Library

📁 **Template**

| Main | Init |
|------|------|
| **Function** | Print |

### Function

**Implemented C functions for CNNs model**
**Currently applied Model : VGG16, VGG19, Resnet 50**
1) Optimized HW function
2) SW function

- MaxPooling2D.txt
- AveragePooling2D.txt
- Convolution2D_Same.txt
- Convolution2D_Valid.txt
  …...
- Dense.txt

# 1. Template Library

📂 **Template**

| Main | Init |
|------|------|
| Function | Print |

**Print**

**Print function
to compare the results of each layer**

- Print4D.txt
- Print3D.txt
- Print2D.txt
- Print1D.txt

# How to match keras and C

```
for (n=0; n<N; n++) {
    for (m=0; m<M; m++) {
        for (x=0; x<F; x++) {
            for (y=0; y<E; y++) {
```

for each output fmap value

convolve
a window
and apply
activation

```
O[n][m][x][y] = B[m];
for (i=0; i<R; i++) {
    for (j=0; j<S; j++) {
        for (k=0; k<C; k++) {
            O[n][m][x][y] += I[n][k][Ux+i][Uy+j] × W[m][k][i][j];
        }
    }
}
```

```
        }
    }
}

O[n][m][x][
        }
    }
}
```

| 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 |

**Keras**

Conv2D(filters, kernel_size, strides, padding, data_format, dilation_rate, activation, use_bias , kernel_initializer, bias_initializer, kernel_regularizer, bias_regularizer, activity_regularizer, kernel_constraint, bias_constraint )
- filters : number of output filter
- kernel_size : filter size
- padding : 'same' , 'valid'
- data_format : 'channels_last' , 'channel_first'
- dilation_rate : k = 1, 2, 4 ..
- activation : activation function
- use bias : W*x + bias ( True/ False )
- kernel_initializer : weight initializer
- ...

# How to match keras and C

## C

```
for (n=0; n<N; n++) {
    for (m=0; m<M; m++) {
        for (x=0; x<F; x++) {
            for (y=0; y<E; y++) {

                max = -Inf;
                for (i=0; i<R; i++) {
                    for (j=0; j<S; j++) {
                        if (I[n][m][Ux+i][Uy+j] > max) {
                            max = I[n][m][Ux+i][Uy+j];
                        }
                    }
                }

                O[n][m][x][y] = max;
            }
        }
    }
}
```

for each pooled value

find the max
with in a window

## Keras

MaxPooling2D(pool-size , strides , padding, data_format )
- pool-size : pooling window
- strides : step
- padding : 'same' ( Add zero-padding )
           'valid'  ( Traverse only with in the image)
- data_format : 'channels_last' (batch, channel, width, height)
               'channels_first' (batch, width, height, channel)

# 2) Result

## 1. Run



```
jiyoung@jiyoung-VirtualBox:~/D

#Run Keras-verification , Generate
#Input : layer information( ex. te
#Output :
#1. Result of C code vs keras code
#2. model_test.cpp , model.cpp

Test_dir="Test_file/Test.csv"
Model_name="Ex_model"
Data_type="ap_uint<16>"
Random_range="20"
return_dir="../"


cd keras-verification
./Verifier.sh $return_dir$Test_dir
cd ../c-code-generation
python Test_cpp_Generator.py $retu
python Cpp_Generator.py $return_di
```

```
#test_dir="Test-file/Test.csv"

Input_file="init_Input.txt"
Weight_file="init_Weight.txt"
Bias_file="init_Bias.txt"

Variable_dir="../Variable_Generator/"
Result_dir="vimdiff.txt"
return_dir="../"

#Generate Variable
#Input: layer info / Output : c_verifier.cpp
cd $Variable_dir
g++ Variable_Generator.cpp -o out
./out $1 $Weight_file $Bias_file $Input_file $2

#Generate C_Verifier
cd ../keras-verification
python C_Verifier_Generator.py $1

#Run Verifier
python Keras_Verifier.py $1 $Variable_dir$Weight_file $Variable_dir$Bias_file $V
ariable_dir$Input_file
g++ -std=c++0x C_Verifier.cpp -o out
./out $Variable_dir$Weight_file $Variable_dir$Bias_file $Variable_dir$Input_file

#Compare result
vimdiff Output/keras_output.txt Output/C_output.txt
```

# 2) Result

## 2. Print Diff

# 2) Result

## 3. Generate File



```cpp
#include <iostream>
#include <ap_int.h>
#include <hls_stream.h>

typedef ap_uint<16> DATA_T;

typedef ap_uint<256> uint256_t;
typedef ap_uint<512> uint512_t;

void Stream_input(DATA_T I[3][10][10], hls::stream<DATA_T> &I_strm) {
  int m, x, y;

#pragma HLS ARRAY_PARTITION variable=I complete dim=1
        Stream_input_x_loop: for (x=0; x<10; x++) {
                Stream_input_y_loop: for (y=0; y<10; y++) {
#pragma HLS PIPELINE

                        Stream_input_m_loop: for (k=0; k<3; k++) {
                                I_strm.write(I[k][x][y]);
                        }
                }
        }
}

void Stream_output(hls::stream<DATA_T> &O_strm, DATA_T O[3]) {
  int x;
#pragma HLS ARRAY_PARTITION variable=O complete dim=1

        Stream_input_x_loop: for (x=0; x<3; x++) {

#pragma HLS PIPELINE

                        O[x] = O_strm.read()
        }
}

static DATA_T I_i[3][10][10];
static DATA_T W1_i[4][3][3][3];
static DATA_T B1_i[4];
static DATA_T W3_i[8][4][3][3];
static DATA_T B3_i[8];
static DATA_T W8_i[8][8];
static DATA_T B8_i[8];
static DATA_T W9_i[3][8];
```

< **Model.cpp** >

```cpp
#include <iostream>
#include <ap_int.h>
#include "hls_stream.h"
#include <stdio.h>
#include <stdlib.h>
using namespace std;

typedef ap_uint<16> DATA_T;

void Ex_model_top(DATA_T I[3][10][10],DATA_T W1[4][3][3][3], DATA_T B1[4],DATA_T
 W3[8][4][3][3], DATA_T B3[8],DATA_T W8[8][8], DATA_T B8[8],DATA_T W9[3][8], DAT
A_T B9[3], DATA_T O[3]);
void Ex_model_sw(DATA_T I[3][10][10],DATA_T W1[4][3][3][3], DATA_T B1[4],DATA_T
W3[8][4][3][3], DATA_T B3[8],DATA_T W8[8][8], DATA_T B8[8],DATA_T W9[3][8], DATA
_T B9[3], DATA_T O[3]);

// argv[1] = init_weight.txt , argv[2] = init_bias.txt , argv[3] = init_input.tx
t
int main(int argc, char *argv[]){

  int m, x, y, i, j, k;
  DATA_T temp;

  static DATA_T I[3][10][10];
  static DATA_T O0_SW[3][10][10];
  static DATA_T W1[4][3][3][3];
  static DATA_T B1[4];
  static DATA_T W3[8][4][3][3];
  static DATA_T B3[8];
  static DATA_T W8[8][8];
  static DATA_T B8[8];
  static DATA_T W9[3][8];
  static DATA_T B9[3];
  static DATA_T O_SW[3];
  static DATA_T O_HW[3];


  FILE *w_stream = fopen(argv[1], "r");
  if (w_stream == NULL) printf("weight file was not opened");
  FILE *b_stream = fopen(argv[2], "r");
  if (b_stream == NULL) printf("bias file was not opened");
  FILE *i_stream = fopen(argv[3], "r");
  if (i_stream == NULL) printf("i_stream file was not opened");
```

< **Testbench.cpp** >

# Conclusion

# Current State

Xiao, Qingcheng, et al. "Exploring heterogeneous algorithms for accelerating deep convolutional neural networks on FPGAs." *Proceedings of the 54th Annual Design Automation Conference 2017*. ACM, 2017.
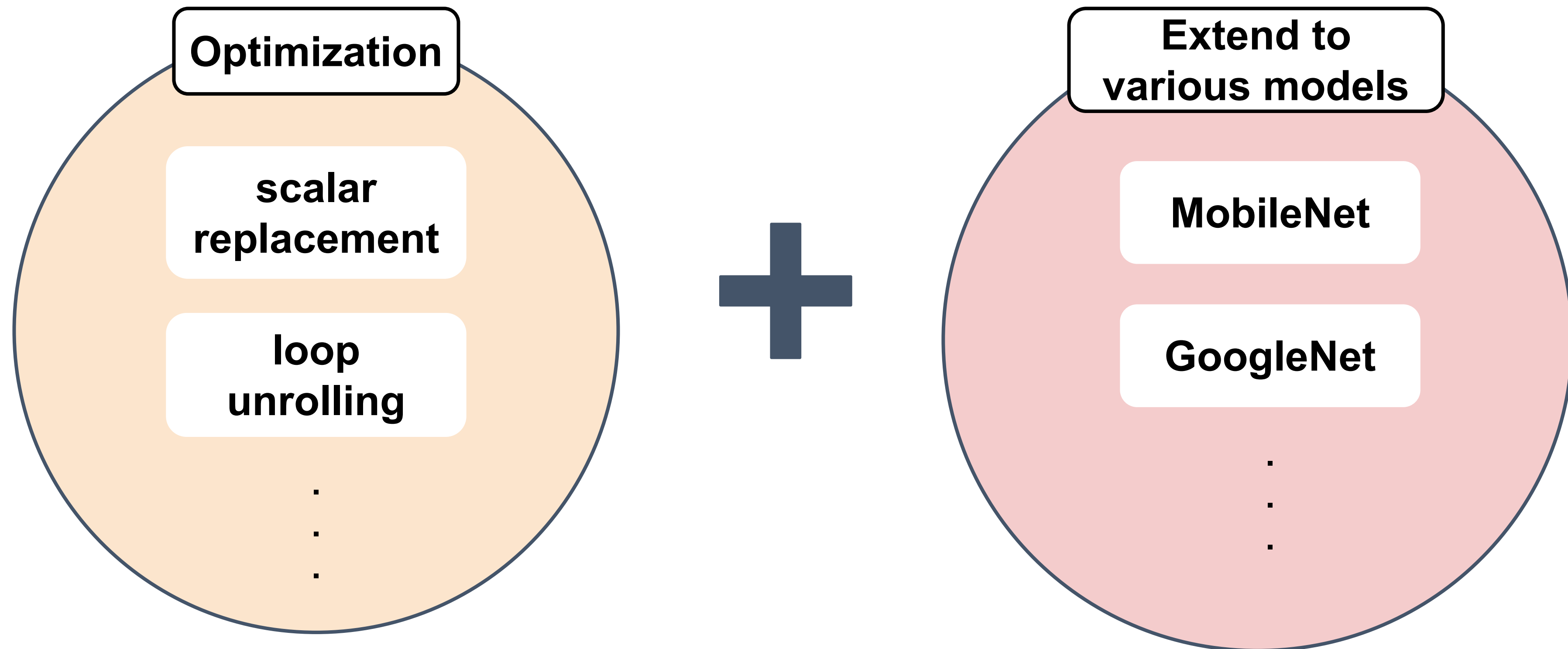
**Reduce memory size**

**Use new optimizing ways**

**Minimize computational latency**

# Future Plan

**Optimization**

scalar
replacement

loop
unrolling

.
.
.

**+**

**Extend to
various models**

MobileNet

GoogleNet

.
.
.

# THANK YOU