# Kneron Inc

## KL520 DME Tutorial

**Kneron Inc**

Kneron

# Table of Contents

# 1. Support New Models in DME Mode

This chapter describes the overview flow in DME mode, and the key points which need to be pay attention to during implementation. The example in this chapter is based on Tiny-YOLO-V3 model, and the dimension of Input Node is 3x224x224 (Channel*Height*Width) in ONNX format.

## 1.1 DME Overview Chart

The below diagram shows the flow chart from host to NPU and the key components involved in DME mode.

### 1.1.1 Host Side

In DME mode, Host is responsible for communicating with SCPU, sending DME commands, firmware information (fw_info.bin), model data (all_models.bin), configuration, and input images to SCPU, and retrieving back the output results from SCPU through USB interface.

### 1.1.2 SCPU

In DME mode, SCPU is responsible for parsing DME commands from Host, writing firmware information, model data, and input images to assigned memory address, launching application, communicating with NCPU, and sending back the output results to Host through USB interface.

### 1.1.3 NCPU and NPU

In DME mode, NCPU is responsible for communicating with SCPU and NPU, registering new preprocess/cpu_op/postprocess functions, setting model, running preprocess/cpu_op/post_process, and enabling NPU. NPU is responsible for running npu_op, interrupting when meeting cpu_op, and continuing running npu_op.

### 1.1.4 Memory

The DDR memory is the space which can be accessed by SCPU, NCPU and NPU. The firmware information, model data, input images, and output results are located in this space.

## 1.2 Limitations for KL520 SDK

(1) The input width and height of models is limited to be less than 256 if keeping aspect ratio.

In this tutorial, the input dimension is 3x224x224 for Tiny-YOLO-V3.

(2) The input width and height of model larger than 255 is supported when doing inference with 1 model, and aspect ratio is changed by default in image preprocessing.

(3) Right shift in image preprocessing is no supported by NPU.

For Tiny-YOLO-V3, the input pixel value range is 0~1 (x/256, x ranges from 0~255). After compiling with compiler, we get that the input_radix is 7 and input_scale is 1.0 for this model. If using NPU to do preprocessing with RGB data, the data is converted to fixed point with the formula x/256 * 2^radix * scale (equals to x/2), which means that each value need right-shift 1 bit. As NPU does not support right-shift for the input data, there are 2 methods to support this model: (a) bypass preprocess and provide the preprocessed RGBA image input directly; (b) do right-shift with NCPU (time-consuming). In this tutorial, method (a) and (b) are introduced for the preprocessing of Tiny-YOLO-V3.

## 1.3 CPU and Output Node Info

The CPU and Output Node info is shown in ioinfo.csv. For Tiny-YOLO-V3, the ioinfo is summarized as below. From this table, cpu_op upsample need to be implemented and there are 2 output nodes used for postprocessing in this model.

| Node Type | Index | Name | Note |
|-----------|-------|------|------|
| c | 0 | up_sampling2d_1_o0 | Upsample CPU Node |
| o | 0 | conv2d_10_o0 | First Output Node |
| o | 1 | conv2d_13_o0 | Second Output Node |

## 1.4 Generate New Configuration Files

Use script dme_conf_generator.py to generate configuration files. Example code is in 2.3.1.

## 1.4.1 Parameter Description

| Parameters | Description |
|---|---|
| Models_selection | Select models to be used in DME |
| output_num | Total output number for this model |
| image_col | Width of the input image |
| image_row | Height of the input image |
| image_ch | Channel of the input image (3 for RGB565, 4 for RGBA) |
| image_format | Input image format which contains settings of "do preprocess" / "bypass preprocess", "output detection results" / "output raw results", "subtract 128 or not", "rotate or not", "parallel processing or not", and input image type. |

## 1.4.2 Image Format

SDK code snippet of Image format flags is defined as below.

```
1.  /* Image format flags */
2.  #define IMAGE_FORMAT_SUB128                BIT31
3.  #define IMAGE_FORMAT_ROT_MASK              (BIT30 | BIT29)
4.  #define IMAGE_FORMAT_ROT_SHIFT            29
5.  #define IMAGE_FORMAT_ROT_CLOCKWISE        0x01
6.  #define IMAGE_FORMAT_ROT_COUNTER_CLOCKWISE 0x02
7.
8.  #define IMAGE_FORMAT_RAW_OUTPUT           BIT28
9.
10. #define IMAGE_FORMAT_CHANGE_ASPECT_RATIO  BIT20
11.
12. #define IMAGE_FORMAT_BYPASS_PRE           BIT19
13. #define IMAGE_FORMAT_BYPASS_NPU_OP        BIT18
14. #define IMAGE_FORMAT_BYPASS_CPU_OP        BIT17
15. #define IMAGE_FORMAT_BYPASS_POST          BIT16
16.
17. #define IMAGE_FORMAT_NPU                  0x00FF
18. #define NPU_FORMAT_RGBA8888       0x00
```

```
19. #define NPU_FORMAT_NIR          0x20
20. /* Support YCBCR (YUV) */
21. #define NPU_FORMAT_YCBCR422     0x30
22. #define NPU_FORMAT_YCBCR444     0x50
23. #define NPU_FORMAT_RGB565       0x60
24.
25. /* Determine the exact format with the data byte sequence in DDR memory: [lowest byte]...[highest byte] */
26. #define NPU_FORMAT_YCBCR422_CRY1CBY0 0x30
27. #define NPU_FORMAT_YCBCR422_CBY1CRY0 0x31
28. #define NPU_FORMAT_YCBCR422_Y1CRY0CB 0x32
29. #define NPU_FORMAT_YCBCR422_Y1CBY0CR 0x33
30. #define NPU_FORMAT_YCBCR422_CRY0CBY1 0x34
31. #define NPU_FORMAT_YCBCR422_CBY0CRY1 0x35
32. #define NPU_FORMAT_YCBCR422_Y0CRY1CB 0x36
33. #define NPU_FORMAT_YCBCR422_Y0CBY1CR 0x37  // Y0CbY1CrY2CbY3Cr...
```

After combining different bit setting, the image_format is generated. The below table shows 8 types of image format which are used in DME mode. 0x10080000, 0x00080000, 0x10000060, and 0x00000060 are used for Tiny-YOLO-V3.

| Image Format | Description | Output Mode |
|---|---|---|
| 0x10080000 | Bypass postprocess, bypass preprocess, RGBA8888 input | Raw outputs |
| 0x00080000 | Bypass preprocess, RGBA8888 input | Detection outputs |
| 0x80000060 | Image pixel value subtract 128, RGB565 input with data byte sequence in DDR memory G0[2:0]B0[4:0]R0[4:0]G0[5:3]G1[2:0]B1[4:0]R1[4:0]G1[5:3]... | Detection outputs |
| 0x80000037 | Image pixel value subtract 128, YCBCR422 (YUV422) input with data byte sequence in DDR memory Y0CbY1CrY2CbY3Cr... | Detection outputs |
| 0x10000060 | Bypass postprocess, RGB565 input | Raw outputs |
| 0x00000060 | RGB565 input | Detection outputs |
| 0x88000060 | Image pixel value subtract 128, parallel processing with 2 image inputs, RGB565 input | Detection outputs |

| 0x98000060 | Image pixel value subtract 128, bypass postprocess, parallel processing with 2 image inputs, RGB565 input | Raw outputs |
|---|---|---|

### 1.4.3 Example Configurations in Host

(1) The below Python code snippet of configuration is for outputting detection results of Tiny-YOLO-V3 with RGBA input. If bypassing preprocess for Tiny-YOLO-V3, the RGBA input is provided.

```
1.  # example for Tiny-YOLO-V3
2.  pre_bypass       = 1
3.  post_bypass      = 0
4.
5.  models_selection = 0x00000001 # Reserved. 0x00000001 if there is only one model
6.  output_num       = 2
7.  image_col        = 224
8.  image_row        = 224
9.  image_ch         = 4
10. image_format     = 0
11. parallel         = 0
12.
13. if pre_bypass:
14.     image_format = image_format | 1 << 19 # 1 << 19 to bypass pre
15.
16. if post_bypass:
17.     image_format = image_format | 1 << 28 # 1 << 28 to bypass post and output raw results
18.
19. if parallel:
20.     image_format = image_format | 1 << 27 # 1 << 27 to do parallel processing
```

(2) The below configuration is for outputting RAW results of Tiny-YOLO-V3 with RGBA input.

```
1.  # example for Tiny-YOLO-V3
2.  pre_bypass       = 1
3.  post_bypass      = 1
4.
5.  models_selection = 0x00000001 # Reserved. 0x00000001 if there is only one model
6.  output_num       = 2
7.  image_col        = 224
```

```
8.  image_row      = 224
9.  image_ch       = 4
10. image_format   = 0
11. parallel       = 0
12.
13. if pre_bypass:
14.     image_format = image_format | 1 << 19 # 1 << 19 to bypass pre
15.
16. if post_bypass:
17.     image_format = image_format | 1 << 28 # 1 << 28 to bypass post and output raw results
18.
19. if parallel:
20.     image_format = image_format | 1 << 27 # 1 << 27 to do parallel processing
```

(3) The below configuration is for outputting detection results of Tiny-YOLO-V3 with RGB565 input.

```
1.  # example for Tiny-YOLO-V3
2.  pre_bypass       = 0
3.  post_bypass      = 0
4.
5.  models_selection = 0x00000001 # Reserved. 0x00000001 if there is only one model
6.  output_num       = 2
7.  image_col        = 640
8.  image_row        = 480
9.  image_ch         = 3
10. image_format     = 0
11. parallel         = 0x00000060
12.
13. if pre_bypass:
14.     image_format = image_format | 1 << 19 # 1 << 19 to bypass pre
15.
16. if post_bypass:
17.     image_format = image_format | 1 << 28 # 1 << 28 to bypass post and output raw results
18.
19. if parallel:
20.     image_format = image_format | 1 << 27 # 1 << 27 to do parallel processing
```

(4) The below configuration is for outputting RAW results of Tiny-YOLO-V3 with RGB565 input.

```
1.  # example for Tiny-YOLO-V3
2.  pre_bypass       = 0
```

```
3.   post_bypass       = 1
4.
5.   models_selection = 0x00000001 # Reserved. 0x00000001 if there is only one model
6.   output_num        = 2
7.   image_col         = 640
8.   image_row         = 480
9.   image_ch          = 3
10.  image_format      = 0
11.  parallel          = 0x00000060
12.
13.  if pre_bypass:
14.      image_format = image_format | 1 << 19 # 1 << 19 to bypass pre
15.
16.  if post_bypass:
17.      image_format = image_format | 1 << 28 # 1 << 28 to bypass post and output raw results
18.
19.  if parallel:
20.      image_format = image_format | 1 << 27 # 1 << 27 to do parallel processing
```

### 1.4.4 Struct to Store the Configuration in SCPU

```
1.   /* KDP image configuration structure */  // SDK code snippet
2.   struct kdp_img_conf_s {
3.       int32_t     image_col;
4.       int32_t     image_row;
5.       int32_t     image_ch;
6.       uint32_t    image_format;
7.       uint32_t    image_memory_address; // not use in DME mode
8.   };
9.
10.  typedef struct kdp_img_conf_s kdp_img_conf_t;
11.
12.  /* KDP dme configuration structure */
13.  struct kdp_dme_conf_s {
14.      int32_t     models_selection;
15.      int32_t     output_num;
16.      struct kdp_img_conf_s img_conf;
17.  };
18.
```

## 1.5 Node Result Layout in Memory

(1) For KL520, the results are in 16-byte aligned format for each row.

The following figure shows that if width mod 16 is equal to 7, then 9 bytes is inserted at the end of each row.



```
1.  // SDK code snippet
2.  /**
3.   * pad_up_16() - calculate the 16-bytes aligned width
4.   *
5.   * @a: the original width
6.   *
7.   * Return value:
8.   * >0 : aligned width
9.   */
10. inline static int pad_up_16(int a)
11. {
12.     return ceil((float)a / 16) * 16;
13. }
```

(2) For KL520, the data sequence of the CPU and Output node is HxCxW (Height*Channel*Width), which means that the data is stored in the sequence: row_1 of C1, row_1 of C2, ..., row_1 of Cn, row_2 of C1, row_2 of C2, ..., row_2 of Cn, row_H of C1, row_H of C2, ..., row_H of Cn.

(4) For KL520, the data type of CPU Node and Output Node is fixed-point (**1 byte for int8_t or 2 bytes for int16_t**).

The below SDK code snippet is to get the data_size.

```
1.  data_size = (POSTPROC_OUTPUT_FORMAT(image_p) & BIT(0)) + 1;     /* 1 or 2 in bytes */
```

To do cpu operation or postprocess, the fixed-point data (**int8_t or int16_t**) need to be converted to float data with do_div_scale() function.

```
1.  // SDK code snippet
2.  /**
3.   * do_div_scale() - convert fixed point to float
4.   *
5.   * @v: fixed point value
6.   *
7.   * @div: 2 raised to the power of radix
8.   *
9.   * @scale: scale for fixed point
10.  *
11.  * Return float value
12.  */
13. static float do_div_scale(float v, int div, float scale)
14. {
15.     return ((v / div) / scale);
```

16. }

For example, the dimension of one output node for Tiny-YOLO-V3 is 7x255x7 (HxCxW), the actual memory size is 7x255x16. The dimension of another output node is 14x255x14 (HxCxW), the actual memory size is 14x255x16. The cpu_op and postprocess functions are based on this data layout.

## 1.6 Register Preprocess/CPU_Operation/Postprocess Functions

### 1.6.1 API Functions

| API Functions | Parameters |
|---|---|
| int kdpio_cpu_op_register(int cpu_op, cpu_op_fn cof); | cpu_op: unique ID for this cpu_node<br><br>cof: CPU operation function, such as **nearest_upsample_cpu** |
| int kdpio_pre_processing_register(int model_id, pre_proc_fn ppf); | model_id: unique ID for this model<br><br>ppf: preprocess function, such as **kdp_preproc_inproc, preprocess_rgba_right_shift_yolo** |
| int kdpio_post_processing_register(int model_id, post_proc_fn ppf); | model_id: unique ID for this model<br><br>ppf: postprocess function, such as **post_yolo_v3** |
| int kdp_preproc_inproc(int model_id, struct kdp_image_s *image_p); | model_id: unique ID for this model<br><br>image_p: kdp_image_s struct with original image input |

kdp_preproc_inproc() supports: Kneron preprocess method (RGB/256 - 0.5), Yolo preprocess method (RGB/256).

### 1.6.2 Allocate Unique ID for New Model and New CPU_OP

(1) Allocate new model ID in /common/include/model_type.h.

In SDK, a new model ID 19 is assigned for TINY-YOLO-V3. This ID is consistent with the model ID in fw_info.bin.

```
1.   enum model_type {
2.       INVALID_ID,
3.       KNERON_FDSMALLBOX              = 1,
4.       KNERON_FDANCHOR               = 2,
5.       KNERON_FDSSD                  = 3,
6.       AVERAGE_POOLING               = 4,
7.       KNERON_LM_5PTS                = 5,
8.       KNERON_LM_68PTS               = 6,
9.       KNERON_LM_150PTS              = 7,
10.      KNERON_FR_RES50               = 8,
11.      KNERON_FR_RES34               = 9,
12.      KNERON_FR_VGG10               = 10,
13.      KNERON_TINY_YOLO_PERSON       = 11,
14.      KNERON_3D_LIVENESS            = 12,
15.      KNERON_GESTURE_RETINANET      = 13,
16.      TINY_YOLO_VOC                 = 14,
17.      IMAGENET_CLASSIFICATION_RES50 = 15,
18.      IMAGENET_CLASSIFICATION_RES34 = 16,
19.      IMAGENET_CLASSIFICATION_INCEPTION_V3= 17,
20.      IMAGENET_CLASSIFICATION_MOBILENET_V2= 18,
21.      TINY_YOLO_V3                  = 19,
22. };
```

(2) Allocate new CPU_OP ID in /ncpu/rtos2/cpu_ops_ex.h.

For a new CPU node, the first step is to get support by Compiler. Then compiler assigns a unique cpu_op_type to this type of CPU node when generating all_models.bin.

As there is a new cpu_op nearest_upsample for TINY-YOLO-V3, cpu_op_type 100 is assigned for it. This op_type is consistent with the one in all_models.bin.

```
1.   /* Type of CPU Operations */
2.   enum cpu_op_type {
3.       SOFTMAX,
```

```
4.      FLATTERN            = 1,
5.      DROPOUT             = 2,
6.      NEAREST_UPSAMPLE    = 100,
7.      BILINEAR_UPSAMPLE   = 101,
8.      MYSTERY             = 1000,
9.  };
```

### 1.6.3 Struct kdp_image_s

The struct kdp_image_s is a key structure which exchange data and parameters between different API functions of NCPU: kdpio_set_model, kdpio_run_preprocess, kdpio_run_npu_op, kdpio_run_cpu_op, and kdpio_run_postprocess.

```
1.  /* Structure of kdp_model_dim */
2.  struct kdp_model_dim_s {
3.      /* CNN input dimensions */
4.      uint32_t    input_row;
5.      uint32_t    input_col;
6.      uint32_t    input_channel;
7.  };
8.
9.  /* Structure of kdp_pre_proc_s */
10. struct kdp_pre_proc_s {
11.     /* input image in memory for NPU */
12.     uint32_t    input_mem_addr;
13.     int32_t     input_mem_len;
14.
15.     /* Input working buffers for NPU */
16.     uint32_t    input_mem_addr2;
17.     int32_t     input_mem_len2;
18.
19.     /* number of bits for input fraction */
20.     uint32_t    input_radix;
21.
22.     /* Other parameters for the model */
23.     void        *params_p;
24. };
25.
26. /* Structure of kdp_post_proc_s */
27. struct kdp_post_proc_s {
28.     /* output number */
29.     uint32_t    output_num;
30.
31.     /* output data memory from NPU */
```

```
32.     uint32_t    output_mem_addr;
33.     int32_t     output_mem_len;
34.
35.     /* result data memory from post processing */
36.     uint32_t    result_mem_addr;
37.     int32_t     result_mem_len;
38.
39.     /* data memory for post processing */
40.     uint32_t    output_mem_addr3;
41.
42.     /* output data format from NPU
43.      *      BIT(0): =0, 8-bits
44.      *              =1, 16-bits
45.      */
46.     uint32_t    output_format;
47.
48.     /* output node parameter */
49.     struct out_node_s    *node_p;
50.
51.     /* Other parameters for the model */
52.     void        *params_p;
53. };
54.
55. /* Structure of kdp_cpu_op_s */
56. struct kdp_cpu_op_s {
57.     /* cpu op node parameter */
58.     struct cpu_node_s    *node_p;
59. };
60.
61. /* KDP image structure */
62. struct kdp_image_s {
63.     /* Original image and model */
64.     struct kdp_img_raw_s   *raw_img_p;
65.     struct kdp_model_s     *model_p;
66.
67.     int        model_id;
68.     /* Setup memory address contains information for input node, cpu node, output node */
69.     char        *setup_mem_p;
70.
71.     /* Model dimension */
72.     struct kdp_model_dim_s  dim;
73.
74.     /* Pre process struct */
75.     struct kdp_pre_proc_s   preproc;
76.
77.     /* Post process struct */
78.     struct kdp_post_proc_s   postproc;
79.
```

```
80.    /* CPU operation struct */
81.    struct kdp_cpu_op_s    cpu_op;
82. };
```

### 1.6.4 Register Preprocess Function

If the new model's preprocess method is not supported by **kdp_preproc_inproc**, the corresponding preprocess function need to be implemented and registered.

The preprocess function can be implemented in /ncpu/rtos2/pre_process_ex.c and registered in /ncpu/rtos2/main.c. The example code is in part (3) of Chapter 2.2.3.

(1) For Tiny-YOLO-V3, preprocess function **kdp_preproc_inproc** is registered for it in /ncpu/rtos2/main.c.

(2) For Tiny-YOLO-V3, another preprocess function **preprocess_rgba_right_shift_yolo** is also implemented in pre_process_ex.c as an example of how to write a new preprocess function. If using this function, the input image is RGBA image with dimensions (HWC, 224x224x4) and each R/G/B value ranges from 0 to 255.

(3) Only 1 preprocess function is registered for each model.

### 1.6.5 Register CPU_Operation Function

By following the results data layout in KL520 and cpu_op algorithm, the corresponding CPU_Operation function can be implemented and registered in SDK.

(1) Nearest Upsample

The following figure shows the nearest upsampling with scale factor of 2. For Tiny-YOLO-V3, the cpu_op function nearest_upsample_cpu() is implemented in /ncpu/rtos2/cpu_ops_ex.c and registered in /ncpu/rtos2/main.c. The example code is in Chapter 2.2.4.

**Nearest**

Input: 2 x 2          Output: 4 x 4

### 1.6.6 Register Postprocess Function

Different models have their specific postprocess method to output the detection results. By following results layout in KL520 and the postprocess algorithm, the corresponding postprocess function can be implemented and registered in SDK.

If bypass postprocess, SDK output the raw results of all output nodes. And no postprocess function is needed to be registered for any models.

For Tiny-YOLO-V3, the postprocess function post_yolo_v3() is implemented in /ncpu/rtos2/post_processing_ex.c and registered in /ncpu/rtos2/main.c. The example code is in Chapter 2.2.5.

### 1.7 Inference with New Models by Calling APIs in Host Library (Serial Mode)

After preparing the input files and code changes to support new model in SDK, call APIs in host library to run test for new model in DME mode.

| Input Files | Description | How to generate |
|---|---|---|
| image.bin | Image input binary file | Convert jpg/png file into binary file (RGB565/RGBA/YCbCr422) |
| all_models.bin | Model data file | Generate by toolchain |

| fw_info.bin | Firmware information data file | Generate by toolchain |
|---|---|---|
| config.bin | Configuration binary file | Use script dme_conf_generator.py to generate |

(1) In order to get the matching results between toolchain and SDK, the preprocessed RGBA image input need to be generated by img_preprocess.py in toolchain. Use 'python img_preprocess.py -h' for more details.

## 1.8 Results to Host

### 1.8.1 Detection Results

The detection results layout follows the definition of struct kdp_app_dme_res_t in kdp_app_dme.h: TOTAL_CLASS (4 bytes, int) + BOXES_COUNT (int) + BOXES_COUNT * BOXES (x1 (float), y1 (float), x2 (float), y2 (float), score (float), class (int)).

```
1.  typedef struct kdp_app_bounding_box_s {
2.      float x1;          // top-left corner:  x
3.      float y1;          // top-left corner:  y
4.      float x2;          // bottom-right corner:  x
5.      float y2;          // bottom-right corner:  y
6.
7.      float score;       // probability score
8.      int class_num;     // class # (of many) with highest probability
9.  } kdp_app_bounding_box_t;
10.
11. #define OBJECT_DETECTION_MAX    80
12.
13. struct kdp_app_dme_res_s {
14.     uint32_t    class_count;            // total class count
15.     uint32_t    box_count;              /* boxes of all classes */
16.     kdp_app_bounding_box_t  boxes[OBJECT_DETECTION_MAX];       /* [box_count] */
17. } kdp_app_dme_res_t;
```

### 1.8.2 RAW Results

The RAW results layout is: OUTPUT_NUM (4 bytes, int) + OUTPUT_NUM * HCW ([height (int), channel, width, h2, c2, w2, ...]) + RAW_DATA (RAW1 (float), RAW2, ...).

## 1.9 Parallel Processing (Parallel Mode)

This chapter focuses on the introduction of parallel processing with 2 input images (RAW images or preprocessed RGBA images) and 2 output memory buffers.

### 1.9.1 Parallel Processing Overview Chart



### 1.9.2 Limitations

(1) Current parallel processing part in parallel processing example SDK is an example of how to do parallel processing and get the timing results. The APIs in host library and logic of sending the postprocessing results to host is not implemented.

(2) Parallel processing in SDK does not always bring the performance improvement due to the presence of CPU operation in preprocessing (such as right-shift), CPU node, time-consuming data movement, etc.

(3) Current example design is to support running parallel processing with 1 model repeatedly.

### 1.9.3 Input and Output Memory

| Memory | To Store | Defined in |
|---|---|---|
| input_mem_addr | Preprocessed RGBA input1 | struct kdp_model_s of ipc.h |
| input_mem_addr2 | Preprocessed RGBA input2 | struct scpu_to_ncpu_s of ipc.h |
| KDP_DDR_BASE_IMAGE_BUF | Raw image input1 | kdp_ddr_table.h |
| KDP_DDR_BASE_IMAGE_BUF2 | Raw image input2 | kdp_ddr_table.h |
| output_mem_addr2 | Previous output data of all output nodes | struct scpu_to_ncpu_s of ipc.h |

Output_mem_addr2 is used as new base address in postprocessing when using parallel processing. The actual addresses of Output nodes are calculated by adding output_mem_addr2 and offsets. The offsets are calculated by subtracting model buffer address from output node address. The below SDK code snippet shows how to copy data of all output nodes to output_mem_addr2.

```
1.  static void npu_out_data_move(void *dst_p)
2.  {
3.      /* Move data out from npu output buffer for parallel handling */
4.      int idx = 0;
5.      int len, data_size, grid_w, grid_h, grid_c, grid_w_bytes_aligned;
6.      int offset = sizeof(struct out_node_s);
7.      int8_t *src_p;
8.      struct out_node_s *out_p;
```

```
9.      data_size = (POSTPROC_OUTPUT_FORMAT(pp_img_p) & BIT(0)) + 1;    /* 1 or 2 in bytes */
10.
11.     for (idx = 0; idx < POSTPROC_OUTPUT_NUM(pp_img_p); idx++) {
12.         out_p = (struct out_node_s *)((uint32_t)POSTPROC_OUT_NODE(pp_img_p) + idx * offset);
13.         if (idx != 0) {
14.             offset = idx * sizeof(struct out_node_s);
15.             out_p = (struct out_node_s *)((uint32_t)POSTPROC_OUT_NODE(pp_img_p) + offset);
16.         } else {
17.             out_p = POSTPROC_OUT_NODE(pp_img_p);
18.         }
19.         src_p = (int8_t *)OUT_NODE_ADDR(out_p);
20.         grid_w = OUT_NODE_COL(out_p);
21.         grid_h = OUT_NODE_ROW(out_p);
22.         grid_c = OUT_NODE_CH(out_p);
23.
24.         len = grid_w * data_size;
25.         grid_w_bytes_aligned = pad_up_16(len);
26.
27.         int total = grid_h*grid_c*grid_w_bytes_aligned;
28.         int mem_offset = OUT_NODE_ADDR(out_p) - MODEL_BUF_ADDR(pp_img_p);  // calculate offset
29.         memcpy((int8_t *)((uint32_t)dst_p + mem_offset), src_p, total);    // copy to new address
30.     }
31.     /* Notify scpu now */
32.     scu_ipc_trigger_to_scpu_int();
33. }
```

### 1.9.4 Struct scpu_to_ncpu_s

This struct contains the global parameters to be passed from SCPU to NCPU when enabling parallel processing: parallel_start, parallel_total, parallel_count, input_mem_addr2, output_mem_addr2, etc.

```
1.  /* Structure of sCPU->nCPU Message */
2.  struct scpu_to_ncpu_s {
3.      ...
4.
5.      /* Images being processed via IPC */
6.      int32_t          active_img_index;
7.      struct kdp_img_raw_s raw_images[IPC_IMAGE_MAX];
8.
9.      /* parallel processing */
10.     uint32_t   parallel_start;  // start parallel processing or not (1 for start, 0 for not start)
11.     uint32_t   parallel_total;  // total times to run parallel processing, 10 in SDK
12.     uint32_t   parallel_count;  // count the finished cycle (from pre to post)
```

```
13.
14.     /* Input/Output working buffers for NPU */
15.     uint32_t    input_mem_addr2; // second RGBA input buffer if using parallel processing
16.     int32_t     input_mem_len2;
17.
18.     uint32_t    output_mem_addr2; // second output buffer is using parallel processing
19.     int32_t     output_mem_len2;
20.
21.     uint32_t    output_mem_addr3;
22. };
```

### 1.9.5 Struct kdp_img_raw_s

(1) This struct contains the parameters to be used in NCPU when enabling parallel processing: parallel_switch_input, parallel_start, parallel_total, parallel_count, tick_start_first_pre, and tick_end_last_post.

(2) The method to generate configuration for parallel processing is shown in chapter 1.4. Bit 27 in image_format is used to enable or disable parallel processing.

| Parameters | Description |
|---|---|
| parallel_switch_input | Parameter for ping-pong operations to switch between input1 and input2 |
| parallel_start | Start parallel processing or not (1 for start, 0 for not start) |
| parallel_total | Total times to run parallel processing (10 in SDK) |
| parallel_count | Count the finished cycle (from preprocessing to postprocessing) |
| tick_start_first_pre | Record the tick for the start of first preprocessing |
| tick_end_last_post | Record the tick for the end of last postprocessing |

```
1.  * Raw image structure */
2.  struct kdp_img_raw_s {
3.      ...
4.
5.      /* Test: NCPU parallel processes */
6.      uint32_t    parallel_switch_input;   // switch to use input1 or input2
7.      uint32_t    parallel_start;   // start parallel processing or not (1 for start, 0 for not start)
8.      uint32_t    parallel_count;   // total times to run parallel processing, 10 in SDK
```

```
9.      uint32_t    parallel_total;   // count the finished cycle (from pre to post)
10.     uint32_t    tick_start_first_pre;    // tick for the start of first preprocessing
11.     uint32_t    tick_end_last_post;      // tick for the end of last postprocessing
12. };
```

### 1.9.6 Code Change for Existing Postprocess Function

The change in postprocess function is the method to calculate the address of output node. This method is compatible with serial mode. The example code is shown as below.

```
1. //src_p        pointer to the output data
2. //image_p      pointer to kdp_image_s
3. //out_p        pointer to out_node_s
4. int8_t
   *src_p = (int8_t *)(POSTPROC_OUTPUT_MEM_ADDR(image_p) + OUT_NODE_ADDR(out_p) - MODEL_BUF_ADDR(image_p));
```

### 1.9.7 Input Files for Parallel Processing

2 input binary files are needed when running parallel preprocessing.

| Input Files | Description | How to generate |
|---|---|---|
| image1.bin and image2.bin | Image input binary files | Convert jpg/png file into binary files (RGB565/RGBA/YCbCr422) |
| all_models.bin | Model data file | Generate by toolchain |
| fw_info.bin | Firmware information data file | Generate by toolchain |
| config.bin | Configuration binary file | Use script dme_conf_generator.py to generate |

## 2. DME Example Code

This chapter focuses on the explanation of example code for Tiny-YOLO-V3. Tiny-YOLO-V3 is a model for multiple-classes object detection. The following sample code implements the object detection function for 80 classes.

## 2.1 SCPU

### 2.1.1 cmd_parser Running in Host Communication Thread

The parser function is to parse different types of commands from Host, write firmware information, model data, and input images to assigned memory address, launch application, and send back the output results to Host through USB interface.

| DME Commands | Description |
|---|---|
| CMD_DME_START | Send START command with FW_INFO size, FW_INFO_DATA (fw_info.bin) and MODEL_DATA (all_models.bin) from Host to KL520 through USB interface |
| CMD_DME_CONFIG | Send CONFIG command with output_num, image_input_w/h/c, and image_format from Host to KL520 through USB interface |
| CMD_DME_SEND_IMAGE | Send IMAGE command with IMAGE_DATA from Host to KL520 through USB interface, and do inference |
| CMD_ACK_NACK | Send ACK command and retrieve back output results from KL520 to Host through USB interface |

(1) CMD_DME_START

After receiving CMD_DME_START from Host, SCPU read firmware information in the message from Host, and write it into the assigned address of DDR memory; then SCPU read the model data from USB interface, and write it into the assigned address of DDR memory: firmware info to KDP_DDR_MODEL_INFO_TEMP, model data to command start address.

```
1.   //Code snippet from cmd_parser() in host_com.c
2.   case CMD_DME_START:
3.   {
4.       // NOTE: Need to determine if DME structure needs to be set up first before processing images
5.       info_msg("[Host_com] Start Object Detection in DME mode.\r\n");
6.       //Get size of inference mode (4 bytes)
7.       uint32_t model_size = msgcmd_arg->addr;
8.
9.       // Load firmware info. The address and size of firmware info is inside the command
10.      memcpy((uint32_t*)KDP_DDR_MODEL_INFO_TEMP, &(msgcmd_arg->len), msghdr->mSize - 4);
11.
12.      uint32_t cmd_addr_offset = 9*sizeof(uint32_t); //cmd addr offset in fw_info
13.      uint32_t cmd_start_addr = *(uint32_t*)(KDP_DDR_MODEL_INFO_TEMP + cmd_addr_offset);
14.
15.      if (model_size <= 0) {
16.          msgrsp_arg->error = MSG_DME_DATA_FAIL;
17.          msgrsp_arg->bytes = CMD_DME_START;
18.          msg_pack(CMD_ACK_NACK, NULL, 0);
19.          break;  // early abort
20.      }
21.
22.      //Load setup, weights and commands
23.      ret = host_mem_read(cmd_start_addr, model_size);  // specify how many bytes
24.      if (ret == model_size) {
25.          APP_id = 3;  // enter DME mode
26.          msgrsp_arg->error = NO_ERROR;
27.          msgrsp_arg->bytes = ret;
28.          msg_pack(CMD_DME_START, NULL, 0);
29.      } else {
30.          msgrsp_arg->error = MSG_DME_DATA_FAIL;
31.          msgrsp_arg->bytes = 0;
32.          msg_pack(CMD_DME_START, NULL, 0);
33.      }
34.
35.      kdp_app_dme_mode(1/*is_dme*/);
36.      break;
37. }
```

## (2) CMD_DME_CONFIG

After receiving CMD_DME_CONFIG from Host, SCPU read configuration information in the message from Host, and copy it to the global variable struct kdp_dme_conf_s.

```
1.   //Code snippet from cmd_parser() in host_com.c
```

```
2.  case CMD_DME_CONFIG:
3.  {
4.      info_msg("[Host_com] DME Config.\r\n");
5.      if (APP_id == 3) {
6.          memcpy(&dme_conf, &(msgcmd_arg->addr), sizeof(struct kdp_dme_conf_s));
7.          msgrsp_arg->error = NO_ERROR;
8.      }
9.      else
10.         msgrsp_arg->error = MSG_APP_NOT_CONFIG;  // not in DME mode
11.     msgrsp_arg->bytes = 0;
12.     msg_pack(CMD_DME_CONFIG, NULL, 0);  // use zero payload response
13.     break;
14. }
```

## (3) CMD_DME_SEND_IMAGE

After receiving CMD_DME_SEND_IMAGE from Host, SCPU read image input from USB interface, and copy it to different assigned address of DDR memory according to configuration: RGBA image input to DDR_MEM_BASE; other types of image input to KDP_DDR_BASE_IMAGE_BUF. Then SCPU allocates memory for final results in advance (the size of allocated memory is greater or equal to the exact size of results), launch kdp_app_dme application to do inference, and send back the exact size of results to Host. Code snippet from cmd_parser() in host_com.c is shown as below.

```
1.  case CMD_DME_SEND_IMAGE:
2.  {
3.      dbg_msg("Received cmd: CMD_DME_SEND_IMAGE\n");
4.
5.      if ((0 == dme_conf.output_num) || (APP_id != 3)) { // check if we have DME structures set up?
6.          if (APP_id == 3)
7.              msgrsp_arg->error = MSG_DME_NO_CONFIG;  // image processing not set up yet
8.          else
9.              msgrsp_arg->error = MSG_APP_NOT_CONFIG;  // not in DME mode
10.         msgrsp_arg->bytes = CMD_DME_SEND_IMAGE;  // add source id
11.         msg_pack(CMD_ACK_NACK, NULL, 0);
12.         break;  // early abort
13.     }
14.
15.     if (!(dme_conf.img_conf.image_format & IMAGE_FORMAT_BYPASS_PRE)) {
16.         // write RGB565 image into input base addr1 of DDR
17.         ret = host_mem_read(KDP_DDR_BASE_IMAGE_BUF, msgcmd_arg->addr);  // specify how many bytes
18.         // write RGB565 image into input base addr2 of DDR if parallel
19.         if (dme_conf.img_conf.image_format & IMAGE_FORMAT_PARALLEL_PROC) {
20.             ret = host_mem_read(KDP_DDR_BASE_IMAGE_BUF2, msgcmd_arg->addr);  // specify how many bytes
```

```
21.          }
22.     } else {
23.          // write RGBA image into input addr1 of DDR
24.          ret = host_mem_read(DDR_MEM_BASE, msgcmd_arg->addr);  // specify how many bytes
25.          // write RGBA image into input addr2 of DDR if parallel
26.          if (dme_conf.img_conf.image_format & IMAGE_FORMAT_PARALLEL_PROC) {
27.               ret = host_mem_read(get_out_com_p()->input_mem_addr2, msgcmd_arg-
     >addr);   // specify how many bytes
28.          }
29.     }
30.
31.     uint32_t raw_len = 0; //For allocating memory in advance
32.     uint32_t output_num = dme_conf.output_num;
33.
34.     // Support parallel processing with 2 inputs, and reset the execution times to 10
35.     if (dme_conf.img_conf.image_format & IMAGE_FORMAT_PARALLEL_PROC) {
36.          get_out_com_p()->parallel_start = 1;
37.          get_out_com_p()->parallel_total = 10;
38.          get_out_com_p()->parallel_count = 0;
39.     } else {
40.          get_out_com_p()->parallel_total = 0;
41.     }
42.
43.     if (dme_conf.img_conf.image_format & IMAGE_FORMAT_RAW_OUTPUT) {
44.          //Get output length in fw_info
45.          uint32_t output_len = *(uint32_t*)(KDP_DDR_MODEL_INFO_TEMP + 6*sizeof(uint32_t));
46.          raw_len = (1 + 3*output_num + output_len)*sizeof(uint32_t);
47.          if (p_raw_out_s == 0) {
48.               p_raw_out_s = kdp_ddr_malloc(raw_len);
49.          }
50.          ret = kdp_app_dme((kdp_app_dme_res_t *)p_raw_out_s, dme_conf);
51.     } else {
52.          if(p_od_out_s == 0) {
53.               p_od_out_s = (kdp_app_dme_res_t *) kdp_ddr_malloc(
54.                    sizeof(kdp_app_dme_res_t));
55.          }
56.          ret = kdp_app_dme(p_od_out_s, dme_conf);
57.     }
58.
59.     //For sending results to HOST
60.     if (dme_conf.img_conf.image_format & IMAGE_FORMAT_RAW_OUTPUT) {
61.          int output_len = 0;
62.          for (int i = 0; i < output_num; i++){
63.               int h = *(uint32_t *)(p_raw_out_s + (3*i+1)*sizeof(uint32_t));
64.               int c = *(uint32_t *)(p_raw_out_s + (3*i+2)*sizeof(uint32_t));
65.               int w = *(uint32_t *)(p_raw_out_s + (3*i+3)*sizeof(uint32_t));
66.               output_len = output_len + h * c * w;
67.          }
68.
```

```
69.        //output_num+HCW+H2C2W2+DATA1+DATA2
70.        final_len = (1 + 3*output_num + output_len)*sizeof(uint32_t);
71.    } else {
72.        int det_num = p_od_out_s->box_count;
73.        final_len = 2*sizeof(uint32_t) + det_num * sizeof(kdp_app_bounding_box_t);
74.    }
75.    msgrsp_arg->error = NO_ERROR;
76.    msgrsp_arg->bytes = final_len;
77.
78.    //Msg for the length of raw results
79.    msg_pack(CMD_DME_SEND_IMAGE, NULL, 0);
80.    break;
81. }
```

## (4) CMD_ACK_NACK

After receiving CMD_ACK_NACK from Host, SCPU send back the results to Host. Depending on the configuration in CMD_DME_CONFIG step, detection results or raw results are sent back to Host.

```
1.  //Code snippet from cmd_parser() in host_com.c
2.  case CMD_ACK_NACK:
3.  {
4.      dbg_msg("[Host_com] Received ACK/NACK\r\n");
5.      // Send results back to host
6.      if (msgcmd_arg->addr) {
7.          no_host_reply();
8.          break;  // if NACK, ignore the command and no reply
9.      }
10.     // We don't enforce Source ID check yet, since only DME is using this
11.     if (dme_conf.img_conf.image_format & IMAGE_FORMAT_RAW_OUTPUT){
12.         data_write((u8 *) p_raw_out_s, final_len);
13.     } else {
14.         data_write((u8 *) p_od_out_s, final_len);
15.     }
16.     break;
17. }
```

## 2.2 NCPU

## 2.2.1 SCPU-to-NCPU Communication Thread

This thread is to register preprocess/cpu_op/postprocess function, wait for all events from SCPU, configure kdp_image_s struct, set model, and run preprocess/npu_op.

(1) Part 1 of scpu_comm_thread

Preprocess/cpu_op/postprocess functions are registered in this part of code.

```
1.  //SCPU Communication and Main Thread
2.  void scpu_comm_thread(void *argument)
3.  {
4.      int32_t flags;
5.      struct kdp_img_raw_s *raw_img_p;
6.      struct kdp_model_s *model_p;
7.      int model_slot_index, image_index;
8.
9.      /* Init kdpio lib and interrupt/status */
10.     kdpio_init();
11.
12.     /* Register preprocess function dynamically */
13.     // Register kdp_preproc_inproc for Tiny-YOLO-V3
14.     kdpio_pre_processing_register(TINY_YOLO_V3, kdp_preproc_inproc);
15.     /* An example to register preprocess function in pre_processing_ex.c */
16.     //kdpio_pre_processing_register(TINY_YOLO_V3, preprocess_rgba_right_shift_yolo);
17.     /* Register cpu_op function dynamically */
18.     kdpio_cpu_op_register(NEAREST_UPSAMPLE, nearest_upsample_cpu);
19.     /* Register post_processing function dynamically */
20.     kdpio_post_processing_register(TINY_YOLO_V3, post_yolo_v3);
21.
22.     // Enable ARM Snoop Control Unit (SCU)'s inter-processor communication to SCPU interrupt
23.     scu_ipc_enable_to_scpu_int();
24.
25.     // Enable SCPU interrupt request
26.     NVIC_EnableIRQ((IRQn_Type)SCPU_IRQ);
27.     // Enable NPU interrupt request
28.     NVIC_EnableIRQ((IRQn_Type)NPU_NPU_IRQ);
29.
30.     // Set NPU to SCPU ID and status
31.     out_comm_p->id = NCPU2SCPU_ID;
32.     out_comm_p->status = STATUS_OK;
33.
```

```
34.    critical_msg("ncpu: Ready!\n");
35.
36.    /* Let's run tickless */
37.    OS_Tick_Disable();
```

## (2) Part 2 of scpu_comm_thread

NCPU wait for all events from SCPU infinitely. If the event flag is
FLAG_SCPU_COMM_ISR, NCPU configures kdp_image_s struct, set model, and run
preprocess/npu_op.

```
1.    /* Waiting for all events from SCPU */
2.    for (;;) {
3.        // Wait forever for Thread Flags of SCPU thread to become signaled
4.        flags = osThreadFlagsWait(FLAG_SCPU_ALL, osFlagsWaitAny, osWaitForever);
5.        // Clear thread flags
6.        osThreadFlagsClear(flags);
7.
8.        // Check the flags to determine whether it is a Interrupt Service Routines from SCPU
9.        if (flags & FLAG_SCPU_COMM_ISR) {
10.            dbg_msg("ncpu: got IPC interrupt\n");
11.
12.            // Check the command from SCPU. If it is CMD_RUN_NPU, execute:
13.            // - Set model and write weights to DMA
14.            // - Preprocess image input
15.            // - Start NPU
16.            if (in_comm_p->cmd == CMD_RUN_NPU) {
17.                // Enable os tick for profiling or disable for saving power
18.                if (log_get_level_ncpu() >= LOG_PROFILE)
19.                    OS_Tick_Enable();
20.                else
21.                    OS_Tick_Disable();
22.
23.                // Get the model index. If there is only 1 model, the index is 0.
24.                model_slot_index = in_comm_p->model_slot_index;
25.                // Get the model
26.                model_p = &in_comm_p->models[model_slot_index];
27.
28.                dbg_msg("NPU: model_slot_index=%d, NPU output mem addr=0x%x, len=%d\n", model_slot_index,
       (uint32_t)model_p->output_mem_addr, model_p->output_mem_len);
29.                dbg_msg("NPU: input_mem_addr = 0x%x, len=0x%x, cmd_mem_addr=0x%x, len=0x%x, wt_addr=0x%x,
       wt_len=0x%x, setup_addr=0x%x, setup_len=0x%x\n",
30.                    (uint32_t)model_p->input_mem_addr, (uint32_t)model_p->input_mem_len,
31.                    (uint32_t)model_p->cmd_mem_addr, (uint32_t)model_p->cmd_mem_len,
32.                    (uint32_t)model_p->weight_mem_addr, (uint32_t)model_p->weight_mem_len,
```

```
33.                  (uint32_t)model_p->setup_mem_addr, (uint32_t)model_p->setup_mem_len);
34.
35.          // Get the image index. If there is only 1 model, the index is 0.
36.          image_index = in_comm_p->active_img_index;
37.          // Get the image input
38.          raw_img_p = &in_comm_p->raw_images[image_index];
39.          npu_img_p = &image_s[npu_img_idx];
40.          npu_img_idx = !npu_img_idx;
41.
42.          npu_img_p->raw_img_p = raw_img_p;
43.          npu_img_p->model_id = in_comm_p->models_type[model_slot_index];
44.          // Set the input memory addr and input memory length for kdp_image_s
45.          // For parallel processing
46.          if (in_comm_p->output_mem_len2 && (RAW_FORMAT(npu_img_p) & IMAGE_FORMAT_BYPASS_PRE)) {
47.              if (0 == RAW_PARALLEL_SWITCH_INPUT(npu_img_p)) {
48.                  PREPROC_INPUT_MEM_ADDR(npu_img_p) = model_p->input_mem_addr;
49.                  PREPROC_INPUT_MEM_LEN(npu_img_p) = model_p->input_mem_len;
50.              } else {
51.                  PREPROC_INPUT_MEM_ADDR(npu_img_p) = in_comm_p->input_mem_addr2;
52.                  PREPROC_INPUT_MEM_LEN(npu_img_p) = in_comm_p->input_mem_len2;
53.              }
54.              RAW_PARALLEL_SWITCH_INPUT(npu_img_p) = !RAW_PARALLEL_SWITCH_INPUT(npu_img_p);
55.          } else {
56.              PREPROC_INPUT_MEM_ADDR(npu_img_p) = model_p->input_mem_addr;
57.              PREPROC_INPUT_MEM_LEN(npu_img_p) = model_p->input_mem_len;
58.          }
59.
60.          // Set the result memory addr and result memory length for kdp_image_s.
61.          // This memory space is used to output result from NPU to NCPU and SCPU
62.          POSTPROC_RESULT_MEM_ADDR(npu_img_p) = raw_img_p->results[model_slot_index].result_mem_addr;
63.          POSTPROC_RESULT_MEM_LEN(npu_img_p) = raw_img_p->results[model_slot_index].result_mem_len;
64.
65.          dbg_msg("NPU: model_type = %d\n",npu_img_p->model_id);
66.
67.          out_comm_p->img_results[image_index].seq_num = raw_img_p->seq_num;
68.          out_comm_p->img_results[image_index].status = IMAGE_STATE_ACTIVE;
69.
70.          // Set the second input memory addr and input memory length for kdp_image_s
71.          // This memory space is used for parallel postprocessing which has not been supported.
72.          PREPROC_INPUT_MEM_ADDR2(npu_img_p) = in_comm_p->input_mem_addr2;
73.          PREPROC_INPUT_MEM_LEN2(npu_img_p) = in_comm_p->input_mem_len2;
74.
75.          // Set the third output memory addr for kdp_image_s
76.          // This memory space is used by postprocess function which need a large space (> 20k bytes) to store
77.          // the parameters.
78.          POSTPROC_OUTPUT_MEM_ADDR3(npu_img_p) = in_comm_p->output_mem_addr3;
```

```
79.
80.              // Set model to be used for processing the input image
81.              kdpio_set_model(npu_img_p, model_p);
82.
83.              // Run preprocess
84.              kdpio_run_preprocess(npu_img_p);
85.              // Run NPU Operation
86.              // - If there is no CPU Node, NPU finishes all NPU Operations and get ready for all Output Node
87.              // - If there is CPU Node, NPU finishes all NPU Operations before the first CPU Node.
88.              kdpio_run_npu_op(npu_img_p);
89.          }
90.
91.          // SCU clear SCPU interrupt
92.          scu_ipc_clear_from_scpu_int();
93.          // Nested Vectored Interrupt Controller (NVIC) clear SCPU interrupt from pending
94.          NVIC_ClearPendingIRQ((IRQn_Type)SCPU_IRQ);
95.          // NVIC enable SCPU interrupt
96.          NVIC_EnableIRQ((IRQn_Type)SCPU_IRQ);
97.      }
98.  }
99. }
```

## 2.2.2 NCPU-to-NPU Communication Thread

This thread is to wait for all events from NPU and run cpu_op/postprocess.

NCPU wait for all events from NPU infinitely. If the event flag is FLAG_NPU_COMM_ISR, NCPU run cpu_op/postprocess, and triggers SCPU interrupt after postprocessing.

```
1.  //NCPU Communication Thread
2.  void npu_comm_thread(void *argument)
3.  {
4.      int32_t flags;
5.      int rc, i;
6.
7.      /* Waiting for all events from NPU */
8.      for (;;) {
9.          // Wait forever for Thread Flags of NPU thread to become signaled
10.         flags = osThreadFlagsWait(FLAG_NPU_COMM_ISR, osFlagsWaitAny, osWaitForever);
11.         osThreadFlagsClear(flags);
12.
```

```
13.          // Check the flags to determine whether it is a Interrupt Service Routines from NPU
14.          if (flags & FLAG_NPU_COMM_ISR) {
15.              // Run CPU Operation
16.              // - If there is no CPU Node, this function return RET_NO_ERROR
17.              // - If there is CPU Node, NCPU finishes the CPU Operation, continue NPU Operations,
18.              //   and return RET_NEXT_NPU.
19.              rc = kdpio_run_cpu_op(npu_img_p);
20.              if (rc == RET_NO_ERROR) {
21.                  i = npu_img_p->raw_img_p->ref_idx;
22.
23.                  dbg_msg("npu out: img[%d] state: %d -> %d\n", i, out_comm_p-
     >img_results[i].status, IMAGE_STATE_NPU_DONE);
24.
25.                  out_comm_p->img_index_done = i;
26.                  out_comm_p->img_results[i].status = IMAGE_STATE_NPU_DONE;
27.                  // Assign the second output memory addr to kdp_image_s
28.                  // This memory space is used for parallel postprocessing which has not been supported.
29.                  POSTPROC_OUTPUT_MEM_ADDR(npu_img_p) = in_comm_p->output_mem_addr2;
30.                  pp_img_p = npu_img_p;
31.
32.                  // NCPU run postprocess
33.                  rc = kdpio_run_postprocess(pp_img_p, npu_out_data_move);
34.
35.                  if (rc == RET_NO_ERROR) {
36.                      i = pp_img_p->raw_img_p->ref_idx;
37.
38.                      dbg_msg("npu out: img[%d] state: %d -> %d\n", i, out_comm_p-
     >img_results[i].status, IMAGE_STATE_DONE);
39.
40.                      out_comm_p->img_index_done = i;
41.                      out_comm_p->img_results[i].status = IMAGE_STATE_DONE;
42.                      // SCU trigger SCPU interrupt
43.                      scu_ipc_trigger_to_scpu_int();
44.                  }
45.              }
46.
47.          NVIC_ClearPendingIRQ((IRQn_Type)NPU_NPU_IRQ);
48.          NVIC_EnableIRQ((IRQn_Type)NPU_NPU_IRQ);
49.      }
50.  }
51. }
```

## 2.2.3 Pre-processing

Preprocess function is registered in scpu_comm_thread(). Kdp_image_s struct provides all information needed for preprocessing: image format, dimension, cropping parameters,

and padding parameters of input raw image; dimension of input node, radix, scale, input memory address, and output memory address.

(1) kdp_preproc_inproc

Kneron provides preprocess API function kdp_preproc_inproc(). The example of how to register it is shown in 2.2.1.

```
 1.  /**
 2.   * kdp_preproc_inproc() - Kneron preprocessing procedure
 3.   *
 4.   * @image_p: pointer to struct kdp_image with buffer of raw image
 5.   *           and dimension, and data for pre/cpu/post processing.
 6.   *
 7.   * @model_id: the model id this function was registered for
 8.   *
 9.   * This is a preprocess function which uses Kneron NPU to accelerate the
10.   * processing and uses Kneron NCPU to do right-
       shift. It can take parameters passed in to do resize, crop, padding with NPU and do right-shift with
       NCPU.
11.   * normalization like -128 and rotation.
12.   *
13.   * Return value:
14.   * 0 : success
15.   * <0 : error
16.   */
17.  int kdp_preproc_inproc(int model_id, struct kdp_image_s *image_p);
```

The preprocessing operations of kdp_preproc_inproc is in the order of crop (original format, such as RGB565/YCBCR422, by NPU), format (original format to RGBA8888, by NPU), resize (RGBA8888, by NPU), left-shift if needed (RGBA8888, by NPU, not implemented in SDK), subtract (RGBA8888, by NPU), pad (RGBA8888, by NPU), right-shift if needed (RGBA8888, by NCPU).

(2) Use kdp_preproc_inproc to do preprocess

For a model using Kneron preprocess (RGB/256 - 0.5), the input pixel value range is -0.5~0.5 (x/256 - 0.5, x ranges from 0~255). After compiling with compiler, we get that the input_radix is 8 and input_scale is 1.0 for this model. If using NPU to do

preprocessing with RGB data, the data is converted to fixed point with the formula **(x/256 - 0.5) * (2^radix) * scale** (equals to x - 128).

For a model using Yolo preprocess (RGB/256), the input pixel value range is 0~1 (x/256, x ranges from 0~255). After compiling with compiler, we get that the input_radix is 7 and input_scale is 1.0 for this model. If using NPU/NCPU to do preprocessing with RGB data, the data is converted to fixed point with the formula **(x/256) * (2^radix) * scale** (equals to x/2).

The following example shows how to convert 640x480x3 RGB565 data to 224x224x4 RGBA data with Kneron preprocess method.

- Generate configuration file with below configuration. This configuration is used for preprocessing RGB565 input.
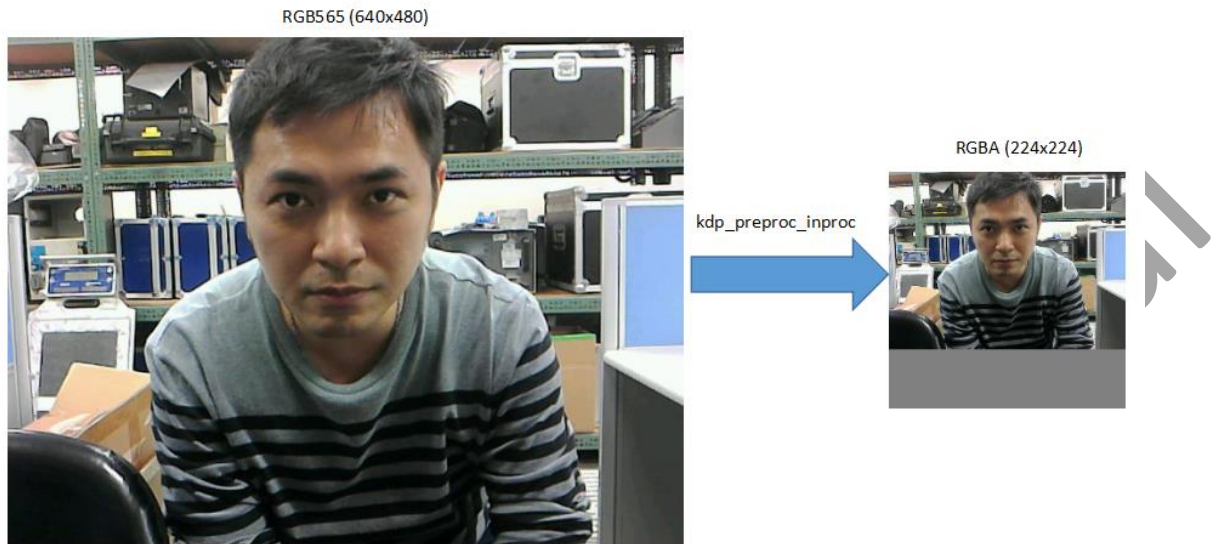
```python
#python code snippet
    output_num  = 1          # (one output for this model)
    image_col   = 640        # (input width)
    image_row   = 480        # (input height)
    image_ch    = 3          # (input channel)
    image_format = 0x80000060 # (bit31 - subtract 128, 0x60 - RGB565 input, not bypass pre/post)
```

- Register preprocess function in scpu_comm_thread() for this model in SDK. The input width and height is 224x224 for MODEL_XXX.

```
/* Register preprocess function dynamically */
kdpio_pre_processing_register(MODEL_XXX, kdp_preproc_inproc);
```

- After sending RGB565 data to KL520 and starting inference, the 640x480 RGB565 data is converted to 224x224 RGBA data by kdp_preproc_inproc.

The preprocessing operations of kdp_preproc_inproc is in the order of crop (top_crop: 0, bottom_crop: 0, left_crop:0, right_crop: 0), format (RGB565 to RGBA8888), resize (640x480 to 224x168 and keep aspect ratio), left-shift (shift: 0), subtract (-128), pad (pad 56 rows with 128), right-shift (shift: 0). The RGBA image below is just for illustration.

RGB565 (640x480)

kdp_preproc_inproc

RGBA (224x224)

(3) preprocess_rgba_right_shift_yolo

preprocess_rgba_right_shift_yolo() is an example of how to implement a customized preprocess function. The example of how to register it is shown in 2.2.1.

```
1.  /**
2.   * preprocess_rgba_right_shift_yolo() - right shift RGBA values for 1 bit
3.   *
4.   * @model_id: the model id this function was registered for
5.   *
6.   * @image_p: pointer to struct kdp_image with buffer of raw image
7.   *           and dimension, and data for pre/cpu/post processing.
8.   *
9.   * This is a preprocess example function which uses Kneron NCPU to do right-
         shift for the input image (HeightxWidthxChannel: 224x224x4).
10.  *
11.  * Return value:
12.  * 0 : success
13.  * <0 : error
14.  */
15.  int preprocess_rgba_right_shift_yolo(int model_id, struct kdp_image_s *image_p)
16.  {
17.      int out_row, out_col;
```

```
18.     int input_radix, bit_shift;
19.     uint8_t *src_p, *dst_p;
20.
21.     //Get model input row and column size
22.     out_row = DIM_INPUT_ROW(image_p);
23.     out_col = DIM_INPUT_COL(image_p);
24.
25.     //Get model input radix
26.     input_radix = PREPROC_INPUT_RADIX(image_p);
27.     //Calculate the number of shift bit
28.     bit_shift = 8 - input_radix; // 1 byte (8 bits) for every R/G/B/A data
29.
30.     //Get source RGBA input address: R/G/B value ranges from 0 to 255
31.     src_p = (uint8_t *)RAW_IMAGE_MEM_ADDR(image_p);
32.     //Get destination RGBA address
33.     dst_p = (uint8_t *)PREPROC_INPUT_MEM_ADDR(image_p);
34.     //Do right-shift
35.     pre_proc_unsign_right_shift(src_p, dst_p, out_row, out_col, bit_shift);
36.
37.     return 0;
38. }
```

(4) Use preprocess_rgba_right_shift_yolo to do preprocess

For a model using Yolo preprocess (RGB/256), the input pixel value range is 0~1 (x/256, x ranges from 0~255). After compiling with compiler, we get that the input_radix is 7 and input_scale is 1.0 for this model. If using NCPU to do preprocessing (only right-shift) with RGBA data (224x224x4), the data is converted to fixed point with the formula **(x/256) * (2^radix) * scale** (equals to x/2).

The following example shows how to convert 224x224x4 RGBA data (R/G/B: 0~255) to 224x224x4 RGBA data  (R/G/B: 0~128).

● Generate configuration file with below configuration. This configuration is used for preprocessing RGBA input (R/G/B: 0~255).

```
1.  #python code snippet
2.      output_num  = 2         # (two outputs for this model)
3.      image_col   = 224       # (input width)
4.      image_row   = 224       # (input height)
5.      image_ch    = 4         # (input channel)
```

```
6.         image_format = 0x00000000 # (0x00 - RGBA input, not bypass pre/post)
```

● Register preprocess function in scpu_comm_thread() for YOLO in SDK. The input width and height is 224x224 for TINY-YOLO-V3. Only one preprocess function is registered for each model.

```
1. /* Register pre_processing function dynamically */
2. //kdpio_pre_processing_register(TINY_YOLO_V3, kdp_preproc_inproc);
3. /* An example to register preprocess function in pre_processing_ex.c */
4. kdpio_pre_processing_register(TINY_YOLO_V3, preprocess_rgba_right_shift_yolo);
```

● After sending RGBA data (R/G/B: 0~255) to KL520 and starting inference, the 224x224x4 RGBA input data is converted to 224x224x4 RGBA data (R/G/B: 0~128) by preprocess_rgba_right_shift_yolo.

### 2.2.4 CPU_OP

CPU_OP function is registered in scpu_comm_thread(). The implementation of this function follows the rule of 16-bytes alignment and HCW data layout. Kdp_image_s struct provides all information needed for cpu_op: dimension of cpu node, radix, scale, input memory address, output memory address.

For Tiny-YOLO-V3, nearest_upsample_cpu() function is implemented and registered. The example of how to register it is shown in 2.2.1.

```
1. /**
2.  * nearest_upsample_cpu() - execute nearest upsampling CPU operations
3.  *
4.  * @image_p: pointer to struct kdp_image with buffer of raw image
5.  *           and dimension, and data for pre/cpu/post processing.
6.  *
7.  * @cpu_op: the CPU operation type id this function was registered for
8.  *
9.  * This is a cpu function which uses Kneron NCPU to do nearest upsampling.
10.  *
11.  * Return value:
12.  * 0 : success
13.  * <0 : error
```

```
14.  */
15.  int nearest_upsample_cpu(int cpu_op, struct kdp_image_s *image_p) {
16.      // Get the input dimension
17.      int row_num_in = CPU_OP_NODE_INPUT_ROW(image_p);
18.      int col_num_in = CPU_OP_NODE_INPUT_COL(image_p);
19.      int chnl_num_in = CPU_OP_NODE_INPUT_CH(image_p);
20.      // Get the output dimension
21.      int row_num_out = CPU_OP_NODE_OUTPUT_ROW(image_p);
22.      int col_num_out = CPU_OP_NODE_OUTPUT_COL(image_p);
23.      // Calculate the input and output 16-bytes aligned column size
24.      int col_num_pad_in = pad_up_16(col_num_in);
25.      int col_num_pad_out = pad_up_16(col_num_out);
26.
27.      int oi, oj, ii, ij, k;
28.      // Calculate the upsampling scales
29.      int h_scale = row_num_out / row_num_in;
30.      int w_scale = col_num_out / col_num_in;
31.      int out_index, in_index;
32.      // Get the input memory address and output memory address
33.      uint8_t *data_array = (uint8_t*)CPU_OP_NODE_INPUT_ADDR(image_p);
34.      uint8_t *dst_data_array = (uint8_t*)CPU_OP_NODE_OUTPUT_ADDR(image_p);
35.      // HCW for DRAM
36.      for (oj = 0; oj < row_num_out; oj++) {
37.          for (k = 0; k < chnl_num_in; ++k) {
38.              for (oi = 0; oi < col_num_out; oi++) {
39.                  ij = oj / h_scale;
40.                  ii = oi / w_scale;
41.                  // Calculate the output index. As each row is 16-bytes aligned, col_num_pad_out is used
42.                  // to calculate the output location
43.                  out_index = oi + col_num_pad_out * (k + chnl_num_in * oj);
44.
45.                  // Calculate the input index. As each row is 16-bytes aligned, col_num_pad_in is used
46.                  // to calculate the input location
47.                  in_index = ii + col_num_pad_in * (k + chnl_num_in * ij);
48.
49.                  // Fill the input data to the output location
50.                  dst_data_array[out_index] = data_array[in_index];
51.              }
52.          }
53.      }
54.      return 0;
55.  }
```

## 2.2.5 Post-processing

Postprocess function is registered in scpu_comm_thread(). The implementation of this function follows the rule of 16-bytes alignment and HCW data layout.

For Tiny-YOLO-V3, post_yolo_v3() function is implemented and registered. The example of how to register it is shown in 2.2.1.

(1) Part 1 of post_yolo_v3

This part of code does preparation for the postprocessing logic: get the result memory address; calculate total class of this model.

```
1.  /**
2.   * post_yolo_v3() - post process function for yolo v3
3.   *
4.   * @model_id: the model id this function was registered for
5.   *
6.   * @image_p: pointer to struct kdp_image with buffer of raw image
7.   *           and dimension, and data for pre/cpu/post processing.
8.   *
9.   * This is a post process function which outputs the detection results.
10.  *
11.  * Return value:
12.  * 0 : success
13.  * <0 : error
14.  */
15. int post_yolo_v3(int model_id, struct kdp_image_s *image_p)
16. {
17.     struct yolo_v3_post_globals_s *gp = &u_globals.yolov3;
18.     int i, j, k, div, ch, row, col, max_score_class, good_box_count, class_good_box_count, good_result_count, len;
19.     float box_x, box_y, box_w, box_h, box_confidence, max_score;
20.     int8_t *src_p, *x_p, *y_p, *width_p, *height_p, *score_p, *class_p;
21.     struct bounding_box_s *bbox;
22.     struct yolo_result_s *result;
23.     struct bounding_box_s *result_box_p, *result_tmp_p, *r_tmp_p;
24.
25.     int data_size, grid_w, grid_h, class_num, grid_w_bytes_aligned, w_bytes_to_skip;
26.     int incr_off;
27.
28.     // Get the data size, the data size is 1 byte for KL520
29.     data_size = (POSTPROC_OUTPUT_FORMAT(image_p) & BIT(0)) + 1;      /* 1 or 2 in bytes */
30.
```

```
31.    // Get the result memory address
32.    result = (struct yolo_result_s *)(POSTPROC_RESULT_MEM_ADDR(image_p));
33.    result_box_p = result->boxes;
34.
35.    result_tmp_p = gp->result_tmp_s;
36.
37.    // Calculate the total classes for this model
38.    class_num = POSTPROC_OUT_NODE_CH(image_p) / YOLO_V3_CELL_BOX_NUM - YOLO_BOX_FIX_CH;
39.
40.    result->class_count = class_num;
41.    bbox = gp->bboxes_v3;
42.    good_box_count = 0;
43.
44.    float anchers_v[3][2];
45.    int idx;
46.    int offset = sizeof(struct out_node_s);
47.    struct out_node_s *out_p;
```

(2) Part 2 of post_yolo_v3

This part of code calculates all valid detection boxes from the 2 output nodes of Tiny-YOLO-V3: calculate the aligned byte size of a row; get the radix and scale of output node; convert fixed point value back to float; calculate all valid detection boxes. Kdp_image_s struct provides all information needed for postprocessing: dimension of output node, radix, scale, output memory address, result memory address.

```
1.    /* Second part of post_yolo_v3 */
2.
3.        // Get the valid detections from the 2 output nodes
4.        for (idx = 0; idx < POSTPROC_OUTPUT_NUM(image_p); idx++) {
5.            // Get the pointer to output node
6.            out_p = (struct out_node_s *)((uint32_t)POSTPROC_OUT_NODE(image_p) + idx * offset);
7.
8.            // Get the pointer to output memory address of output node
9.            src_p = (int8_t *)OUT_NODE_ADDR(out_p);
10.
11.            // Get width and height of output node
12.            grid_w = OUT_NODE_COL(out_p);
13.            grid_h = OUT_NODE_ROW(out_p);
14.            // Calculate the byte size of a row
15.            len = grid_w * data_size;
16.            // Calculate the aligned byte size of a row: 16-byte alignment
17.            grid_w_bytes_aligned = round_up(len);
18.            // Calculate the skipped bytes for each row
```

```
19.        w_bytes_to_skip = grid_w_bytes_aligned - len;
20.        len = grid_w_bytes_aligned;
21.
22.        // Get the start address of x, y, w, h, confidence_score, class probability
23.        x_p = src_p;
24.        y_p = x_p + len;
25.        width_p = y_p + len;
26.        height_p = width_p + len;
27.        score_p = height_p + len;
28.        class_p = score_p + len;
29.
30.        // Get the radix of output node and calculate 2 raised to the power of radix
31.        div = 1 << OUT_NODE_RADIX(out_p);
32.        // Get the scale of output node
33.        float fScale = *(float *)&OUT_NODE_SCALE(out_p);
34.
35.        // Set different anchors for output node
36.        if (0 == idx) {
37.            memcpy(anchers_v, anchers_v0, 6*sizeof(float));
38.        } else {
39.            memcpy(anchers_v, anchers_v1, 6*sizeof(float));
40.        }
41.
42.        // The data layout is in HCW
43.        for (row = 0; row < grid_h; row++) {
44.            for (ch = 0; ch < YOLO_V3_CELL_BOX_NUM; ch++) {
45.                for (col = 0; col < grid_w; col++) {
46.                    if (data_size == 1) {
47.                        box_x = (float)*x_p;
48.                        box_y = (float)*y_p;
49.                        box_w = (float)*width_p;
50.                        box_h = (float)*height_p;
51.                        box_confidence = (float)*score_p;
52.                    } else {
53.                        box_x = (float)*(uint16_t *)x_p;
54.                        box_y = (float)*(uint16_t *)y_p;
55.                        box_w = (float)*(uint16_t *)width_p;
56.                        box_h = (float)*(uint16_t *)height_p;
57.                        box_confidence = (float)*(uint16_t *)score_p;
58.                    }
59.
60.                    // Calculate the confidence score with sigmoid function
61.                    box_confidence = sigmoid(do_div_scale(box_confidence, div, fScale));
62.
63.                    /* Get scores of all class */
64.                    for (i = 0; i < class_num; i++) {
65.                        if (data_size == 1)
66.                            gp->box_class_probs[i] = (float)*(class_p + i * len);
```

```
67.                    else
68.                        gp->box_class_probs[i] = (float)*(uint16_t *)(class_p + i * len );
69.                    // Calculate the class probability with sigmoid function
70.                    gp->box_class_probs[i] = sigmoid(do_div_scale(gp-
    >box_class_probs[i], div, fScale));
71.                }
72.
73.                //increase pointer to next position
74.                x_p += data_size;
75.                y_p += data_size;
76.                width_p += data_size;
77.                height_p += data_size;
78.                score_p += data_size;
79.                class_p += data_size;
80.
81.                /* Find highest score class */
82.                max_score = 0;
83.                max_score_class = -1;
84.                for (i = 0; i < class_num; i++) {
85.                    gp->box_class_probs[i] = gp->box_class_probs[i] * box_confidence;
86.                    if (gp->box_class_probs[i] >= prob_thresh_yolov3) {
87.                        if (gp->box_class_probs[i] >= prob_thresh_yolov3 && gp-
    >box_class_probs[i] > max_score) {
88.                            max_score = gp->box_class_probs[i];
89.                            max_score_class = i;
90.                        }
91.                    }
92.                }
93.                /* Calculate the valid detection box: top-left, right-bottom */
94.                if (max_score_class != -1) {
95.                    box_x = do_div_scale(box_x, div, fScale);
96.                    box_y = do_div_scale(box_y, div, fScale);
97.                    box_w = do_div_scale(box_w, div, fScale);
98.                    box_h = do_div_scale(box_h, div, fScale);
99.
100.                   box_x = (sigmoid(box_x) + col) / grid_w;
101.                   box_y = (sigmoid(box_y) + row) / grid_h;
102.                   box_w = expf(box_w) * anchers_v[ch][0] / DIM_INPUT_COL(image_p);
103.                   box_h = expf(box_h) * anchers_v[ch][1] / DIM_INPUT_ROW(image_p);
104.
105.                   bbox->x1 = (box_x - (box_w / 2)) * DIM_INPUT_COL(image_p);
106.                   bbox->y1 = (box_y - (box_h / 2)) * DIM_INPUT_ROW(image_p);
107.                   bbox->x2 = (box_x + (box_w / 2)) * DIM_INPUT_COL(image_p);
108.                   bbox->y2 = (box_y + (box_h / 2)) * DIM_INPUT_ROW(image_p);
109.                   bbox->score = gp->box_class_probs[max_score_class];
110.                   bbox->class_num = max_score_class;
111.
112.                   bbox++;
113.                   good_box_count++;
```

```
114.                    }
115.                }
116.
117.                /* go to next grid_w*85 values for next anchors*/
118.                // The output dimension is grid_h * 255 * grid_w, and there are 3 anchors
119.                // The data sequence is:
120.                // 1st anchor: row of x, row of y, row of w, row of h, row of confidence score, row of
       class_prob1, ..., row of class_prob80;
121.                // 2nd anchor: row of x, row of y, row of w, row of h, row of score, row of class1, ..
       ., row of class80;
122.                // 3rd anchor: row of x, row of y, row of w, row of h, row of score, row of class1, ..
       ., row of class80;
123.                incr_off = w_bytes_to_skip + grid_w_bytes_aligned*84; //84 = 85 -1
124.                x_p += incr_off;
125.                y_p += incr_off;
126.                width_p += incr_off;
127.                height_p += incr_off;
128.                score_p += incr_off;
129.                class_p += incr_off;
130.            }
131.        }
132.
133.     }
134.
135.     good_result_count = 0;
```

## (3) Part 3 of post_yolo_v3

This part of code does non max suppression (NMS) on all valid detection boxes. Then final detection boxes are generated.

```
1.  /* Third part of post_yolo_v3 */
2.
3.      // Do NMS for all valid detection boxes, and get the final good result
4.      for (i = 0; i < class_num; i++) {
5.          bbox = gp->bboxes_v3;
6.          class_good_box_count = 0;
7.          r_tmp_p = result_tmp_p;
8.
9.          for (j = 0; j < good_box_count; j++) {
10.             if (bbox->class_num == i) {
11.                 memcpy(r_tmp_p, bbox, sizeof(struct bounding_box_s));
12.                 r_tmp_p++;
13.                 class_good_box_count++;
14.             }
```

```
15.           bbox++;
16.       }
17.
18.       if (class_good_box_count == 1) {
19.           memcpy(&result_box_p[good_result_count], &result_tmp_p[0], sizeof(struct bounding_box_s));
20.           good_result_count++;
21.       } else if (class_good_box_count >= 2) {
22.           qsort(result_tmp_p, class_good_box_count, sizeof(struct bounding_box_s), box_score_comparator);
23.           for (j = 0; j < class_good_box_count; j++) {
24.               if (result_tmp_p[j].score == 0)
25.                   continue;
26.               for (k = j + 1; k < class_good_box_count; k++) {
27.                   if (box_iou(&result_tmp_p[j], &result_tmp_p[k]) > nms_thresh_yolov3) {
28.                       result_tmp_p[k].score = 0;
29.                   }
30.               }
31.           }
32.
33.           for (j = 0; j < class_good_box_count; j++) {
34.               if (result_tmp_p[j].score > 0) {
35.                   memcpy(&result_box_p[good_result_count], &result_tmp_p[j], sizeof(struct bounding_box_s));
36.                   good_result_count++;
37.               }
38.           }
39.       }
40.   }
41.
42.   dbg_msg("good_result_count: %d\n", good_result_count);
43.   result->box_count = good_result_count;
44.   for (unsigned int i = 0; i < good_result_count; i++) {
45.       dbg_msg("post_yolo3 %f %f %f %f %f %d\n", result->boxes[i].x1, result->boxes[i].y1, result->boxes[i].x2, result->boxes[i].y2, result->boxes[i].score, result->boxes[i].class_num);
46.   }
47.   len = good_result_count * sizeof(struct bounding_box_s);
48.
49.   return len;
50. }
```

## 2.3 dme_conf_generator.py

The python script is modified accordingly to the model input dimension by referring to 1.4. It generates different binary configuration file which is sent to KL520 through USB interface. And each binary file is consistent with one type of configuration.

```
1.  # scripts to generate configuration file in DME mode
2.
3.  # example for Tiny-YOLO-V3
4.  pre_bypass      = 1
5.  post_bypass     = 0
6.
7.  models_selection = 0x00000001 # Reserved. 0x00000001 if there is only one model
8.  output_num      = 2
9.  image_col       = 224
10. image_row       = 224
11. image_ch        = 4
12. image_format    = 0
13. parallel        = 0
14.
15. if pre_bypass:
16.     image_format = image_format | 1 << 19 # 1 << 19 to bypass pre
17.
18. if post_bypass:
19.     image_format = image_format | 1 << 28 # 1 << 28 to bypass post and output raw results
20.
21. if parallel:
22.     image_format = image_format | 1 << 27 # 1 << 27 to do parallel processing
23.
24. models_selection   = models_selection.to_bytes(4, byteorder='little')
25. output_num         = output_num.to_bytes(4, byteorder='little')
26. image_col          = image_col.to_bytes(4, byteorder='little')
27. image_row          = image_row.to_bytes(4, byteorder='little')
28. image_ch           = image_ch.to_bytes(4, byteorder='little')
29. image_format       = image_format.to_bytes(4, byteorder='little')
30.
31. img_conf = models_selection + output_num + image_col + image_row + image_ch + image_format
32.
33. with open("config.bin", 'wb') as binary_file:
34.     binary_file.write(img_conf)
```