# HiAI DDK V320

# FAQs

**Issue**       03

**Date**      2020-02-28

# Huawei HiAI Application

Send an application email to developer@huawei.com.

Email subject: HUAWEI HiAI + Company name + Product name

Email body: Cooperation company + Contact person + Phone number + Email address

We will reply to you within 5 working days.

Official website: https://developer.huawei.com/consumer/en/

# Change History

Changes between document issues are cumulative. The latest document issue contains all the changes made in earlier issues.

| Date | Version | Change Description |
|---|---|---|
| 2020-02-28 | 03 | Added the description of the .so library compilation mode. |
| 2019-12-31 | 02 | Added the description of HiAI DDK V320. |
| 2019-09-04 | 01 | Added the description of HiAI DDK V310. |

# Contents

# 1 HiAI DDK FAQs

## 1.1 How Do I Check Whether My Smartphone Supports the NPU?

Call the **GetVersion**, **CheckModelCompatibility**, and **Load** methods of class AiModelMngerClient and the **BuildModel** method of class AiModelBuilder as follows.

If running **inference_npu_demo** fails due to a model incompatibility error, inlcude **libcpucl.so** to the demo to run the model on the CPU. For details, see **Android.mk** in **inference_cpu_demo**.

# 1.2 How Do I Check the Time Consumption of Forward Computation?

- In sync mode, add timestamp print before and after the process API at the JNI layer. For details, see **app_sample\inference_demo\Demo_Soure_Code\app\src\main\jni\classify_jni.cpp**.

```
341        // before process
342        struct timeval tpstart, tpend;
343        gettimeofday(&tpstart, nullptr);
344        int istamp;
345        int ret = g_clientSync->Process(context, input_tensor, output_tensor, 1000, istamp);
346        if (ret) {
347            LOGE("[HIAI_DEMO_SYNC] Runmodel Failed!, ret=%d\n", ret);
348            return nullptr;
349        }
350
351        // after process
352        gettimeofday(&tpend, nullptr);
353        float time_use = 1000000 * (tpend.tv_sec - tpstart.tv_sec) + tpend.tv_usec - tpstart.tv_usec;
354
355        time_use_sync =  time_use / 1000;
```

- In async mode, add timestamp print to the callback function and the OnProcessDone API. If the callback is called only once, the time difference between the callback function and the OnProcessDone API is calculated. If the callback is called more than once, the time difference between the callbacks is calculated. For details, see **app_sample\inference_demo\Demo_Soure_Code\app\src\main\jni\classify_async_jni.cpp**.

```
55   void JNIListener::OnProcessDone(const AiContext &context, int result1, const vector<shared_ptr<AiTensor>> &output_tensor, int32_t istamp)
56   {
57       std::unique_lock<std::mutex> lock(mutex_map);
58       map_input_tensor.erase(istamp);
59       condition_.notify_all();
60
61       gettimeofday(&tpend, nullptr);
62       time_use = 1000000 * (tpend.tv_sec - tpstart.tv_sec) + tpend.tv_usec - tpstart.tv_usec;
63       LOGI("[HIAI_DEMO_ASYNC] AYSNC infrence time %f ms.", time_use / 1000);
64       LOGE("[HIAI_DEMO_ASYNC] AYSNC JNI layer onRunDone istamp: %d", istamp);
65
66       JNIEnv *env = nullptr;
```

# 1.3 How Do I Select the Sync or Async APIs?

Both are supported, depending on the service requirements. However, async APIs are recommended for better performance. (Theoretically, if a callback function is to be called, async APIs should be used.)

# 1.4 How the Input Image Format Is Converted Before Model Inference?

During inference, source images are stored in ARGB format.

In non-AIPP scenarios, the images are laid out according to NCHW before being input to the model.

In AIPP scenarios, the images are laid out according to NCHW and then converted into YUV images before being input to the model.

For details about the format conversion code, see **app_sample\inference_demo\Demo_Soure_Code\app\src\main\java\com\huawei\hiaidemo\view\NpuClassifyActivity.java**.

```java
231         @Override
232    ●↑   protected void onActivityResult(int requestCode, int resultCode, Intent data) {
233             super.onActivityResult(requestCode, resultCode, data);
234             if (resultCode == RESULT_OK && data != null) switch (requestCode) {
235                 case GALLERY_REQUEST_CODE:
236                     try {
237                         Bitmap bitmap;
238                         ContentResolver resolver = getContentResolver();
239                         Uri originalUri = data.getData();
240                         bitmap = MediaStore.Images.Media.getBitmap(resolver, originalUri);
241                         String[] proj = {MediaStore.Images.Media.DATA};
242                         Cursor cursor = managedQuery(originalUri, proj,  selection: null,  selectionArgs: null,  sortOrder: null);
243                         cursor.moveToFirst();
244                         Bitmap rgba = bitmap.copy(Bitmap.Config.ARGB_8888,  isMutable: true);
245                         initClassifiedImg = Bitmap.createScaledBitmap(rgba, selectedModel.getInput_W(), selectedModel.getInput_H(),  filter: true);
246                         byte[] inputData = {};
247                         if(selectedModel.getUseAIPP()){
248                             inputData = Untils.getPixelsAIPP(selectedModel.getFramework(), initClassifiedImg, selectedModel.getInput_W(), selectedModel.getInput_H())
249                         }else {
250                             inputData = Untils.getPixels(selectedModel.getFramework(), initClassifiedImg, selectedModel.getInput_W(), selectedModel.getInput_H());
251                         }
252                         ArrayList<byte[]> inputDataList = new ArrayList<>();
253                         inputDataList.add(inputData);
254                         Log.d(TAG,  msg: "inputData.length  1 is :"+inputData.length+"");
255                         runModel(selectedModel, inputDataList);
256                     } catch (IOException e) {
257                         Log.e(TAG, e.toString());
258                     }
```

For details, see the **Utils.java** file in the same directory.

```
74        public static byte[] getPixels(String framework, Bitmap bitmap,
75                                      int resizedWidth, int resizedHeight) {
76            int channel = 1;
77            float[] buff = new float[channel * resizedWidth * resizedHeight];
78
79            int rIndex;
80            int gIndex;
81            int bIndex;
82    /*...*/
97            int pixCount = channel * resizedWidth * resizedHeight;
98            byte[] ret = new byte[pixCount * 4];
99
00            for (int i = 0; i < pixCount; ++i) {
01                int int_bits = Float.floatToIntBits(buff[i]);
02                ret[(i * 4) + 0] = (byte) int_bits;
03                ret[(i * 4) + 1] = (byte) (int_bits >>> 8);
04                ret[(i * 4) + 2] = (byte) (int_bits >>> 16);
05                ret[(i * 4) + 3] = (byte) (int_bits >>> 24);
06            }
07
08            return ret;
09        }


111        public static byte[] getPixelsAIPP(String framework,Bitmap bitmap, int resizedWidth, int resizedHeight){
112            Log.i(TAG,  msg: "resizedWidth : " + resizedWidth +  " resizedHeight : " + resizedHeight);
113            return getNV12(resizedWidth,resizedHeight, bitmap);
114        }
115
116    @   private static byte [] getNV12(int inputWidth, int inputHeight, Bitmap scaled) {
117            // Reference (Variation) : https://gist.github.com/wobbals/5725412
118
119            int [] argb = new int[inputWidth * inputHeight];
120
121            Log.i(TAG,  msg: "scaled : " + scaled);
122            scaled.getPixels(argb,  offset: 0, inputWidth,  x: 0,  y: 0, inputWidth, inputHeight);
123
124            byte [] yuv = new byte[inputWidth*inputHeight*3/2];
125            encodeYUV420SP(yuv, argb, inputWidth, inputHeight);
126
127            //scaled.recycle();
128
129            return yuv;
130        }
131
132        private static void encodeYUV420SP(byte[] yuv420sp, int[] argb, int width, int height) {
```

# 1.5 What Are the Data Formats Supported by the Inference Function?

The data format of **input_tensor** configured in the **process** API must be NCHW.

# 1.6 How Do I Compile .so Libraries?

The DDK HiAI external APIs are encapsulated in the dynamic library **libhiai*.so**.

The standard library **c++_shared** is a dynamic dependency (**APP_STL := c++ _shared**). To avoid problems caused by version differences, you are advised to compile the standard library as a dynamic library.

# 1.7 What Do I Do If Type of the Operator Output by the OMG Offline Model Is Incorrect?

In the Caffe network, some layers of the same type serve different computing purposes. For example, at the DetectionOutput layer, you need to explicitly specify the detection operator types such as FSRDetectionOutput and SSDDetectionOutput through operator mapping. Otherwise, the OMG fails to generate the offline model. You can include the **--op_name_map** parameter to the OMG command. For details, see section "General Parameters" in *Huawei HiAI DDK V320 OMG Tool Instructions*. Alternatively, you can explicitly specify the output operator type in the .proto model file of the original network, for example, SSDDetectionOutput, as shown in the following figure.

**Figure 1-1** Output operator type before and after modification

```
layer {
  name: "detection_out"
  type: "DetectionOutput"
  bottom: "mbox_loc"
  bottom: "mbox_conf_flatten"
  bottom: "mbox_priorbox"
  top: "detection_out"
  include {
    phase: TEST
  }
  detection_output_param {
    num_classes: 21
    share_location: true
    background_label_id: 0
    nms_param {
      nms_threshold: 0.45
      top_k: 400
    }
    save_output_param {
      label_map_file: "data/VOC0712/labelmap_voc.prototxt"
    }
    code_type: CENTER_SIZE
    keep_top_k: 200
    confidence_threshold: 0.3
  }
}
```

```
layer {
  name: "detection_out"
  type: "SSDDetectionOutput"
  bottom: "mbox_loc"
  bottom: "mbox_conf_flatten"
  bottom: "mbox_priorbox"
  top: "detection_out"
  include {
    phase: TEST
  }
  detection_output_param {
    num_classes: 21
    share_location: true
    background_label_id: 0
    nms_param {
      nms_threshold: 0.45
      top_k: 400
    }
    save_output_param {
      label_map_file: "data/VOC0712/labelmap_voc.prototxt"
    }
    code_type: CENTER_SIZE
    keep_top_k: 200
    confidence_threshold: 0.3
  }
}
```