# Ingenic®

## Magik

Post-Training-Quantization User Guide

Date: Feb. 2022

北京君正集成电路股份有限公司
Ingenic Semiconductor Co., Ltd.

# Ingenic®

## Magik Post-Training-Quantization User Guide

**Release history**

| Date | Author | Revision | Change |
|------|--------|----------|--------|
| Feb. 2022 | Owen | 1.0.1 | First release |
| | | | |

## Disclaimer

**Ingenic Semiconductor Co., Ltd.**

目录

# 1 Introduction

The MagikTransformKit tool supports post-quantization of the onnx/tensorflow format model to obtain the magik format model, which can be deployed on the chip side.

The document is based on the yolov5s network to show how to generate the magik format model deployed on the chip, and how to locate the problem when the generated model effect is not ideal.

# 2 Conversion commands

```
1  ./magik-transform-tools \
2  --framework onnx \
3  --target_device T40 \
4  --outputpath yolov5s_magik.mk.h \
5  --inputpath ./yolov5s.onnx \
6  --mean 0,0,0 \
7  --var 255.0,255.0,255.0 \
8  --img_width 640 \
9  --img_height 640 \
10 --img_channel 3 \
11 --input_nodes images \
12 --output_nodes 640,660,680 \
13 --post_training_quantization true \
14 --ptq_config_path ./ptq_yolov5_config.json
```

The MagikTransformKit tool completes the post-trainning-quantization work of the onnx or tensorflow model through command-line parameters, and supports training and quantization model conversion work. The command-line parameters are divided into general parts and post-trainning-quantization parts.

## 2.1 common command line parameters

**framework** – Input model format, now supports onnx/tensorflow format model

(The caffe/mxnet/pytorch format model can uniformly convert the onnx format model)

**target_device** – part number of the deployment board

**inputpath** – Input model to be converted

**outputpath** – The output of the conversion. After the conversion is completed, the required

*.bin file will be generated in the specified directory

**mean** –The mean of the preprocessing at training time

**var** – The variance of the preprocessing at training time

**img_width** – The width of the input when running on the chip side

**img_height** – The height of the input when running on the chip side

**img_channel** – The channel of the input when running on the chip side

**input_nodes** – Model input node

**output_node** – Model output node

## 2.2 Post-training-quantization command line arguments

**post_training_quantization** – The post-training-quantization switch options

**ptq_config_path** – "json" file path that is provided to the post-training-quantization part.

# 3 Json file

The json file consists of two parts, the basic information required by the MagikTransformKit tool for post-training-quantization and the debug information of the debugging model.

When the quantization model generated by the MagikTransformKit tool is not satisfactory, you can add debug information to the Json file for debugging.

## 3.1 Basic information of json

```
1  {
2      "format":"RGB",
3      "width":640,
4      "height":640,
5      "path":"./yolov5-50",
6      "used_image_num":2,
7      "feature_quantize_method":"KL",
8      "weight_quantize_method":"MAX_ABS"
9  }
```

**format** – image processing format: RGB,BGR,GRAY

**width** – Input width of the network inference

**height** – Input height of the network inference

**path** – Directory path of the calibration image set.

**used_image_num** – Number of the calibration image set.

**feature_quantize_method** – feature quantization method: KL, ADMM, MSE

**weight_quantize_method** – weight quantization method: MAX,ADMM

The MagikTransformKit tool gets the correct command line information and json file, and can generate a quantized model with a suffix *.mk.h in magik format and a quantized model with a suffix *.bin, which can be deployed on the board.

```
1  [          _g      yolov5s_demo]$ls
2  magik_model_yolov5s_magik.mk.h   run.sh          yolov5s.onnx
3  magik-transform-tools            yolov5-20
4  ptq_yolov5_config.json           yolov5s_magik.bin
5  [szhang@102 160 1 202 yolov5 demo]$
```

## 3.2 json debug option (optional)

```
1  {
2      "format":"RGB",
3      "width":640,
4      "height":640,
5      "path":"./yolov5-50",
6      "used_image_num":2,
7      "feature_quantize_method":"KL",
8      "weight_quantize_method":"MAX_ABS",
9      "feature_quantize_bits":8,
10     "weight_quantize_bits":8,
11     "debug": {
12         "open":false,
13         "debug_input_path":"./test_image/5050_1-112.jpg",
14         "debug_output_dir":"./dump_data"
15     },
16     "bias_corr":{
17         "weight_corr":false,
18         "feature_corr":false,
19         "vali_dir":"./validated_images",
20         "corr_dir":"./test_image"
21     },
22     "weight_fake_quant": true,
23     "feature_fake_quant": true,
24     "skip_node_quant":[
25     ],
26     "skip_op_quant":[
27     ]
28  }
```

**feature_quantize_bits** – feature quantization bit width, optional, default 8 bits

**weight_quantize_bits** – feature quantization bit width, optional, default 8 bits

**debug** – Option to locate the poor inference effect of the quantitative model, optional{

    **debug_input_path** – path of the debug calibration image set

    **debug_output_dir** – path of the debug results

}

**weight_fake_quant** – Simulate weight for quantization effect, optional

**feature_fake_quant** – Simulate feature for quantization effect, optional

**bias_corr** – Adjust the bias information of the network, optional {

    **weight_corr** – quantize the weight in order to adjust the bias information of the network

    **feature_corr** – quantize the feature in order to adjust the bias information of the network

    **vali_dir** – path of the validation image set

    **corr_dir** – path of the correction bias image set

}

**skip_node_quant** – Skip specifying operator name quantization, optional

**skip_op_quant** – Skip specifying operator type quantization, optional

# 4 yolov5 post-training-quantization demo

Demo path：*magik-toolkit/Models/post/soc_sample/T40/yolov5s/*

## 4.1 Prepare the model

The MagikTransformKit tool directly supports tensorflow/onnx format models for post-quantization work, and indirectly supports Pytorch/mxnet/caffe format models. These three formats need to be uniformly converted to onnx format.

For the caffe format model, we provide the caffe2onnx tool for conversion work (see caffe2onnx_tools for details)

```
1 (test) [          :        caffe2onnx_tools]$python :
  scaffe2onnx/caffe2onnx.py model/resnet18.prototxt model:
  s/resnet18.caffemodel -o resnet18.onnx
```

caffe2onnx path: *magik-toolkit/Models/post/caffe2onnx_tools*

The currently used yolov5s model can be downloaded from Google network disk:
Url:

https://drive.google.com/file/d/1AKZy5rqkq5kalJz6tUqNg_kLV2uf43Nt/view?usp=sh

aring

model download to path: ***magik-toolkit/Models/post/yolov5s/***


## 4.2 Prepare the calibration picture

The calibration pictures in the post-training-quantization process are usually extracted from the training set or the validation set, and the number is between 10 and 100.

Make sure that there are only picture files in the calibration picture directory, and no other files, otherwise the quantization process will go wrong.



## 4.3 Json file

1）Ensure that the format/weight/height parameters are consistent with the training time

2 ） Specify the path of the quantized pictures (19 in total) and the number of quantized pictures 12

3）Commonly used quantitative(feature:KL,weight:MAX_ABS)

```
{
    "format":"RGB",
    "width":640,
    "height":640,
    "path":"./yolov5-20",
    "used_image_num":12,
    "feature_quantize_method":"KL",
    "weight_quantize_method":"MAX_ABS"
}
```

## 4.4 Check command line writing

1）Check that the output model parameter suffix is *.mk.h

2）Set the post_training_quantization option to true, and confirm that the input file path of the

ptq_config_path option is correct.

```
1 ./magik-transform-tools \
2 --framework onnx \
3 --target_device T40 \
4 --outputpath yolov5s magik.mk.h \
5 --inputpath ./yolov5s.onnx \
6 --mean 0,0,0 \
7 --var 255.0,255.0,255.0 \
8 --img_width 640 \
9 --img_height 640 \
10 --img_channel 3 \
11 --input_nodes images \
12 --output_nodes 640,660,680 \
13 --post_training_quantization true \
14 --ptq_config_path ./ptq_yolov5_config.json
```

3）The quantization process of the yolov5 network is performed normally by the MagikTransformKit tool, and the yolov5s_magik.bin file running on the chip will be generated.

```
1 [                        yolov5s_demo]$ls
2 magik_model_yolov5s_magik.mk.h  run.sh        yolov5s.onnx
3 magik-transform-tools          yolov5-20
4 ptq_yolov5_config.json         yolov5s_magik.bin
```

**yolov5s.onnx** ─ Model file to be quantified

**yolov5-20** ─ post-training-quantization using calibration image set

**ptq_yolov5_config.json** ─ post-training-quantization using json file

**run.sh** ─ Script of convert commands

**yolov5s_magik.bin** ─ run file on chip

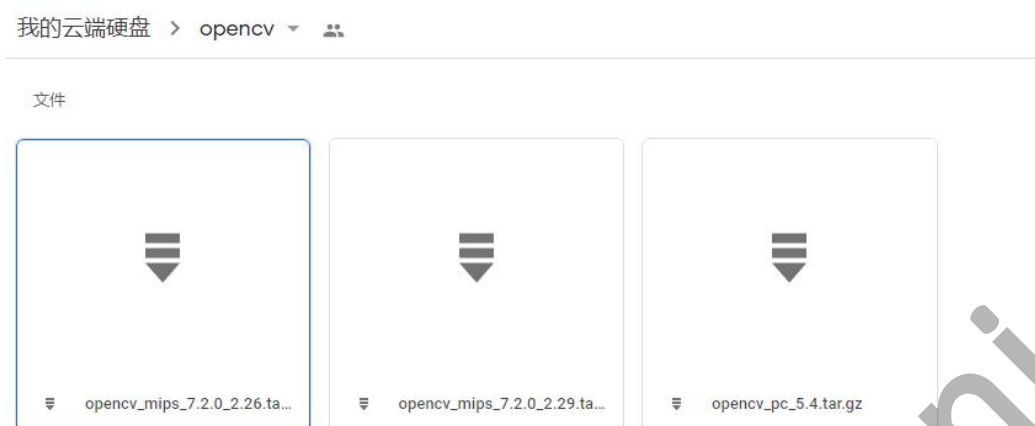**magik_model_yolov5s_magik.mk.h** ─ run file on chip

# 5 Run on-board

## 5.1 Code compling

### 5.1.1 Preparation

First, download the resources compiled by our mips
Link:https://drive.google.com/drive/folders/1Gdj4vF2TpaPzbPNT0H1hIEKDTXEj92Jg?

usp=sharing



Please download and decompress it into magik-toolkit/ThirdParty.

$ **tar -xvf opencv_mips_7.2.0_2.29.tar.gz**

Secondly, we also provide inference.cpp, Makefile and test data, which are necessary while running the model on board. In addition, we also need venus library and mips compilation tool. The venus library is located in: "***magik-toolkit/InferenceKit/nna1/***", and the mips compilation tool is provided during debugging process.

### 5.1.2 Input of network

Here, in order to demonstrate the process quickly, we've added opencv library functions here to facilitate calling "imread()" and other functions, so we can test   jpg diagram directly on board.

### 5.1.3 Model loading

The model in inference.cpp is passed in through parameters. So, please pay attention to synchronously copying it to the corresponding directory on the board side and passing it in at run time.

### 5.1.4 Hyperparameter settings

```
void generateBBox(std::vector<venus::Tensor> out_res, std::vector<magik::venus::ObjBbox_t>& candidate_boxes,
int img_w, int img_h)
{
  float person_threshold = 0.3;
  int classes = 80;
  float nms_threshold = 0.6;
  std::vector<float> strides = {8.0, 16.0, 32.0};
  int box_num = 3;
  std::vector<float> anchor = {10,13,  16,30,  33,23, 30,61,  62,45,  59,119, 116,90,  156,198,  373,326};
}
```

These settings should be consistent with the test code on the PC side.

### 5.1.5 Compile

TOPDIR -- directory of venus

opencv_path -- path of the compiled opencv library

Libtype -- muclibc or other(glibc)

build_type -- release (Default setting.)

        -- profile (Visualization of network structure at runtime and statistics of running time and GOPs at each layer of the network.)

        -- debug    (Save the results of quantization features of each layer.)

        -- nmem (Count the memory usage of nmem when the model is running, and it's in /tmp/nmem_memory.txt.)

Directly compile inference.cpp, then it will generate venus_ yolov5s_ bin_ uclibc_*，which is the executable file we need on the board.

We also provide the code inference_nv12.cpp and input data 10_w1024_ h714. nv12 when you want to use nv12 type input. If necessary, just modify Makefile for compilation and use test.

## 5.2 Run

Opencv library path: *magik-toolkit/ThirdParty/opencv_mips_7.2.0_2.29/uclibc/lib*
Venus library path: *magik-toolkit/InferenceKit/nna1/mips720-glibc229/lib/uclibc/*

### 5.2.1   Release(default)

$ **make build_type=release**

(use "**make build_type=release clean**" if you want to clean)

We can get executable file venus_yolov5s_bin_uclibc_release after the compiling process, then copy libvenus.so, venus_yolov5s_bin_uclibc_release, opencv libraries (which consists of libopencv_core.so.3.3,          libopencv_highgui.so.3.3,          libopencv_imgcodecs.so.3.3, libopencv_imgproc.so.3.3, libopencv_videoio.so.3.3), yolov5s_magik.bin and fall_1054_sys.jpg to board.

When you access the board, please add the library path first (**export LD_LIBRARY_PATH=$lib_path:$LD_LIBRARY_PATH**), after that , you just need to run it on board: "**./venus_yolov5s_bin_uclibc_release yolov5s_magik.bin  fall_1054_sys.jpg**". The operation results are shown in the figure below

```
1 [root@Ingenic-uc1_1:venus_sample_ptq_yolov5s]# ./venus_yolov5s_bin_ptq_uclibc_release .
  ./yolov5s_magik.bin fall_1054_sys.jpg
2 ./venus_yolov5s_bin_ptq_uclibc_re
3 lease ../yolov5s_magik.bin fall_1054_sys.jpg
4 Warning : The version number is not obtained. Please upgrade the soc-nna!
5 INFO(magik): venus version:0.9.0(00000900_8450a71)  built:20220211-1206(7.2.0 glibc2.26
  mips@aie)
6 w:388 h:574
7 ori_image w,h: 388 ,574
8 model-->640 ,640 4
9 input shape:
10 -->384 640
11 scale---> 0.668990
12 resize padding over:
13 resize valid_dst, w:260 h 384
14 padding info top :0 bottom 0 left:190 right:190
15 pad_x:190 pad_y:0 scale:0.668990
16 box:    7 41 352 406 0.88
17 box:   97 329 380 510 0.71
```

### 5.2.2 Debug（Check data）

See step 7.2 of the "Magik Training Quantization User Guide" for details, the processing of the input is a little different.

### 5.2.3 Profile (network visualization)

$ **make build_type=profile**

(use "**make build_type=profile clean**" if you want to clean)

We can get executable file ***venus_yolov5s_bin_uclibc_profile*** after compiling process, then copy libvenus.p.so, venus_yolov5s_bin_uclibc_profile, opencv library(which consists of libopencv_core.so.3.3, libopencv_highgui.so.3.3, libopencv_imgcodecs.so.3.3, libopencv_imgproc.so.3.3, libopencv_videoio.so.3.3), yolov5s_magik.bin and fall_1054_sys.jpg to board.

When you access the board, please add library path first (**export LD_LIBRARY_PATH=$lib_path:$LD_LIBRARY_PATH**), then run the program on the board with: "**./venus_yolov5s_bin_uclibc_profile yolov5s_magik.bin fall_1054_sys.jpg**". The operation results are shown in the figure below.

### 5.2.4   Nmem(nmem occupancy)

$ **make build_type=nmem**

(use "**make build_type=nmemclean**" if you want to clean)

We can get executable file ***venus_yolov5s_bin_uclibc_release*** after compiling process, then copy libvenus.m.so, venus_yolov5s_bin_uclibc_nmem, opencv library(which consists of libopencv_core.so.3.3, libopencv_highgui.so.3.3, libopencv_imgcodecs.so.3.3, libopencv_imgproc.so.3.3, libopencv_videoio.so.3.3), yolov5s_magik.bin and fall_1054_sys.jpg to board.

When you access the board, please add library path first (**export LD_LIBRARY_PATH=$lib_path:$LD_LIBRARY_PATH**), then run the program on the board with: "**./venus_yolov5s_bin_uclibc_nmem yolov5s_magik.bin bus.jpg**". The operation results are shown in the figure below.



# 6  Check data

To verify whether the data on the PC side and the chip side can be aligned, follow the steps below:

## 6.1 On PC

Use the PC-side reasoning library provided by Ingenic, by setting the relevant environment variables. The operator layer-by-layer data can be saved as a binary file and compared with the binary file after the SOC has opened the corresponding mode. The specific operation steps on the PC side are as follows:

In order to keep the input data highly consistent, we use the opencv library as an image parsing tool. If you need to compare the data consistency between python and pc, we recommend that you also use opencv as a tool for parsing images on the python side. Therefore, we have sorted out the necessary libraries and models for running PC-side tools and placed them on Google Cloud.

First, you need to switch to the ***magik-toolkit/Models/post/pc_inference/*** folder, and set the

environment variable on the PC side

**export MAGIK_CPP_DUMPDATA=true**

The quantization results of each layer can be saved to the ./MAGIK_DATA_DUMP/ folder. The specific running command is:

**export   MAGIK_CPP_DUMPDATA=true**

**make    clean**

**make   -j8**

**./pc_inference_bin   ./model_quant.mgk   fall_1054_sys.jpg   640 384**

when the conversion tool converts the model, *fall_1054_sys.jpg* is the input image, and 640 and 384 correspond to the input width and height respectively. Finally save *./MAGIK_DATA_DUMP/* (see the figure below for details), it can be seen that the input layer is saved as *input_data_shape_1_384_640_3.bin*, and the following shape naming rules are n, h, w, c, that is, the height is 384, the width is 640, and the last is the output layer; if you don't want to do it step by step, you can run our provided *extractor.sh* script to run it in one click.



## 6.2 On chip

(1) The input data required by the board side is 4-channel data format. For the "fall_1054_sys.jpg" tested on the PC side above, we use the input "input_data_shape_1_384_640_3.bin" saved at runtime to process it into .h and use it for the board side test to ensure that the inputs of the two

are exactly the same , you can do offline processing first:

Script in:

magik-toolkit//Models/post/soc_sample/T40/yolov5s/venus_sample_ptq_yolov5s/generate_img_input.py:

**python generate_img_input.py bin_path w h c bin**

**bin_path** -- The bin file to be converted, that is, the input bin saved by the PC side operation

**w, h** --The width and height of the bin file correspond to the bin file name, which are 1,h,w,c

**bin** -- The flag sign indicates that the header file is turned from the bin file, and the other is the img header file.

Note: If you use bin, you can completely ensure that the input on both sides is consistent, which is convenient for data.

In the provided example, the model is passed in through parameters (argv[1]). When running, pay attention to synchronously copying it to the corresponding directory on the board, and passing it in as needed at run time. The width (w) and height (h) are also passed in from the outside, and the specific values are the same as those in (1).

Compile c code: **make build_type=debug**

The executable file, "venus_yolov5s_bin_ptq_uclibc_debug" will be generated in the current folder, then copy the venus library (libvenus.d.so), the executable file (venus_yolov5s_bin_ptq_uclibc_debug), model file (yolov5s_magik.bin) to the development board, finally run:

**./venus_yolov5s_bin_ptq_uclibc_debug  yolov5s_magik.bin 640 383**

Among them, 640 corresponds to w, 383 corresponds to h, and the output and other information of each layer will be automatically saved during operation (as shown in the following figure)



Among them, **_*_QuantizeFeature_out.bin** and PC-side runtime are saved in **_*_QuantizeFeature.bin** under **./MAGIK_DATA_DUMP/**

There is a one-to-one correspondence, and you can directly check whether the md5 value is consistent. In the case of ensuring that the input is completely consistent, there are layers in the middle that do not correspond to timely feedback.

# 7 Frequently Asked Questions

Q: The conversion process prompts invalid json

```
2022-02-11 17:15:16.480355: E interface.cc:69
s-parase_input_parater] Invalid json
```

Answer: Check if the json file path and command are correct

```
1  ./magik-transform-tools \
2  --framework onnx \
3  --target_device T40 \
4  --outputpath yolov5s_magik.mk.h \
5  --inputpath ./yolov5s.onnx \
6  --mean 0,0,0 \
7  --var 255.0,255.0,255.0 \
8  --img_width 640 \
9  --img_height 640 \
10 --img_channel 3 \
11 --input_nodes images \
12 --output_nodes output \
13 --post_training_quantization true \
14 --ptq_config_path ./ptq_yolov5_config.jso
15 > > > > > > > > > > > > > > INFO(magik): magik-
   stransform-tools version:1.0.0(00010000_fbdaa2
   s8)  cuda version:9.0.176  built:20220209-1705
   s_GPU
16 2022-02-11 17:15:09.308386: W clip_model_by_n
   sop_optimizer.cc:48-collect_remove_nodes] Mult
   siple output NOp exists on the current pathway
   s!
17 2022-02-11 17:15:16.480355: E interface.cc:69
   s-parase_input_paramter] Invalid json
```

, and check if the json format is correct

```
1 {
2     "format":"RGB",
3     "width":640,
4     "height":640,
5     "path":"./yolov5-50",
6     "used_image_num":2,
7     "feature_quantize_method":"KL",
8     "weight_quantize_method":"MAX_ABS"
9 }
```
✗

```
1 {
2     "format":"RGB",
3     "width":640,
4     "height":640,
5     "path":"./yolov5-50",
6     "used_image_num":2,
7     "feature_quantize_method":"KL",
8     "weight_quantize_method":"MAX_ABS"
9 }
```
✓

# 8 json debug process

Turn on the debug option, MagikTransformKit will simulate the network native model (onnx or Tensorflow) inference and chip-side inference, save the inference results of each layer of operators, and calculate the errors of each layer of operators (cosine similarity, Euclidean distance, absolute maximum value).

After the network is quantized, the accuracy loss is within the acceptable range of 0.05% to 1%. Usually, the cosine similarity performance is as follows:

A. The cosine similarity of the first layer to the last layer is 99.0%+;

B. Or there is a little 98.0%+ in the middle layer, and the cosine similarity of the first and last layers is 99.0%+.

## 8.1 The phenomenon of poor effect

Using the json debug option of MagikTransformKit to locate the problem, the phenomenon can be divided into the following four types:

A. The first layer operator cosine_avg does not reach 99.0%+

B. The quantization of the first layer reaches 99.0%+, a certain layer begins to decline, and the cosine similarity of the last layer is lower than 95.0%

C. The quantization of the first and middle layers is normal, and the quantization cosine similarity of the last few layers decreases

D. The cosine similarity performance is normal, and the chip-side verification effect is poor

## 8.2 Solution

For the above summary phenomenon information, we organize the steps and methods to locate and solve the problem.

A. The first layer operator cosine_avg does not reach 99.0%+ or max_absolu is very large:

Reason: There is a problem with the network input. Under normal circumstances, the cosine_avg of the first layer operator is greater than 99%

Solution: Check the image format (BRG, RGB) and mean variance in the json file to ensure that these parameters are consistent with training

B. The quantization of the first layer reaches 99.0%+, a certain layer begins to decline, and the cosine similarity of the last layer is lower than 95.0%:

Reason 1: Quantization brings errors, which must exist

Solution: Use the feature_quantize_bit or weight_quantize_bit option to increase the bit width (set to 10, 12, 14, default 8 bits)

Reason 2: There is a problem with the quantization method of a certain layer of operators

Solution: Use the skip_node_quant option to skip the serious layer of cosine_avg drop, after using the mechanism, the cosine similarity of other layers is normal

"skip_node_quant":[
    "node_name_0"
],

Feed the problem back to the MagikTransformKit docking staff, and at the same time send the dump results of the input and output of the problematic layer under the path "result/path/ptq_model" and the path "result/path/float_model", and send the quantization parameters of the problem layer to the docking staff

Reason 3: Problems simulating native network inference

MagikTransformKit has problems simulating network native model (onnx or tf) inference. Use the "./result/path/input_uint8.bin" data generated by the conversion tool as the input of the native model inference, and save the results of the native model inference; the generated results and the "result/path/float_model" path store the simulation results (float binary form NHWC ) to compare and confirm whether there is a problem in simulating the native model inference;

If there is a problem and report it to the docking staff, send the input and output bin files of the problem layer and the correct output file under the path "result/path/float_model" to the docking staff to fix the problem.

C. The cosine similarity of the first layer and the middle layer is normal, and the quantized cosine similarity of the last few layers decreases:

Reason 1: Refer to phenomenon 2)

Reason 2: There is a cumulative error in network quantization, and the error caused by layer-by-layer quantization is obvious in a certain layer

Solution: Confirm the error caused by weight quantization or feature quantization through the "weight_fake_quant" option and the "feature_fake_quant" option, and use the "bias_corr" option to eliminate the error

```
"feature_fake_quant": true,
"weight_fake_quant": true,
"bias_corr":{
    "weight_corr":false,
    "feature_corr":false,
    "vali_dir":"./validated_images",
    "corr_dir":"./test_image"
}
```

"weight_corr" option: Eliminate errors caused by weight quantization

"feature_corr" option： Eliminate errors caused by feature quantization

"corr_dir" option: Pass in the image set for bias calibration

"vali_dir" option: Pass in the verified image set, and after the bias calibration is completed, verify the bias calibration effect

Using the "bias_corr" option still cannot solve the problem, and the quantization error can be eliminated by increasing the quantization bit width option.

```
"feature_quantize_bits":10,
"weight_quantize_bits":10,
```

D. The operator cosine similarity is normal, but the accuracy is poor:

Reason 1: There is a problem with the conversion tool simulating native network reasoning,

which is the same as the reason 3 of phenomenon 2), but the manifestation of the phenomenon is different

Solution: Refer to phenomenon 2) Reason 3

Reason 2: There is a problem with inference on the chip side

Solution: Use the "./result/path/input_uint8.bin" data generated by the conversion tool as the input of the onboard inference, and save the result information of the onboard inference output. Compare the generated results with the simulation results stored in the **"result/path/ptq_model"** path (uint8_t or uint16_t binary form NHWC) to confirm whether there is a problem with the on-board inference;

If there is a problem and feedback the conversion tool to the docking staff, send the "***./result/path/input_uint8.bin***" data generated by MagikTransformKit and the dump result of the input and output of the problematic layer under the path "result/path/ptq_model" to the docking staff.