



Kneron Linux Toolchain Manual

2022 Sep Toolchain v0.18.2

PDF Downloads ([../res/manual.pdf](#))

0. Overview

KDP toolchain is a set of software which provide inputs and simulate the operation in the hardware KDP 520, KDP 720 and KDP 530. For better environment compatibility, we provide a docker which include all the dependencies as well as the toolchain software.

This document is compatible with `kneron/toolchain:v0.18.2`.

Performance simulation result on NPU KDP520:

Model	Size	FPS (npu only)	Time(npu only)	Has CPU node(s)?
Inception v3	224x224	26.27	160 ms	No
Inception v4	299x299	1.45	687 ms	No
Mobilenet v1	224x224	57.3	17.4 ms	No
Mobilenet v2	224x224	54.7	18.3 ms	No
Mobilenet v2 ssdlite	300x300	28.5	35.1 ms	No
Resnet50 v1.5	224x224	6.94	144 ms	No
OpenPose	256x256	6.37	1569 ms	No
SRCNN	384x384	11.0	90.9 ms	No
Tiny YOLOv3	416x416	21.8	45.9 ms	Yes
YOLOv3	416x416	1.44	693 ms	Yes
YOLOv5s	640x640	3.67	272 ms	Yes
Lite-HRNet	256x192	8.82	113 ms	Yes

Performance simulation result on NPU KDP720:

Model	Size	FPS (npu only)	Time(npu only)	Has CPU node(s)?
Inception v3	224x224	86.2	11.6 ms	No
Inception v4	299x299	20.4	49.0 ms	No
Mobilenet v1	224x224	437	2.29 ms	No
Mobilenet v2	224x224	677	1.48 ms	No
Mobilenet v2 ssdlite	300x300	310	3.22 ms	No
Resnet50 v1.5	224x224	55.6	18.0 ms	No
OpenPose	256x256	5.30	187 ms	No
SRCNN	384x384	134	7.48 ms	No
Tiny YOLOv3	416x416	151	6.61 ms	No
YOLOv3	416x416	10.1	98.6 ms	No
YOLOv5s	640x640	25.7	38.9 ms	No
Centernet res101	512x512	2.84	352 ms	No
Lite-HRNet	256x192	136	7.38 ms	No

Performance simulation result on NPU KDP530:

Model	Size	FPS (npu only)	Time(npu only)	Has CPU node(s)?
Inception v3	224x224	64.3	15.5 ms	No
Inception v4	299x299	16.5	60.5 ms	No
Mobilenet v1	224x224	289	3.46 ms	No
Mobilenet v2	224x224	340	2.94 ms	No
Mobilenet v2 ssdlite	300x300	205	4.88 ms	No
Resnet50 v1.5	224x224	35.7	28.0 ms	No
OpenPose	256x256	3.61	277 ms	No
SRCNN	384x384	54.5	18.3 ms	No
Tiny YOLOv3	416x416	72.0	13.9 ms	No
YOLOv3	416x416	5.93	169 ms	No
YOLOv5s	640x640	16.9	59.3 ms	No
Centernet res101	512x512	2.19	457 ms	No
Unet	384x384	0.950	1050 ms	No
Lite-HRNet	256x192	252.7	19.0 ms	No

In this document, you'll learn:

1. How to install and use the toolchain docker.
2. What tools are in the toolchain.
3. How to utilize the tools through Python API.

Major changes of the current version

• [v0.18.2]

- ktc: Add `mode`, `optimize`, `export_dynasty_dump` argument to analysis.
- ktc: Set `skip_verify` in analysis as deprecated.
- regression: Add `optimize` option for optimization level selection.
- regression: Fix interface to assure platform is integer.
- converter: Add 720 batch process with `--opt-720` flag.
- converter: Add enable shared weight duplication flag `-d`. By default, shared weights are no longer duplicated.
- converter: Remove `-s` flag since it is now the default behaviour.
- converter: Optimize debug output.

- compiler: Fix RDMA not correctly executed.
- E2E simulator: Change dynasty library fetching method.
- Minor bug fixes.
- **[v0.18.1]**
 - docker: Update numpy from 1.18.5 to 1.21.
 - ktc: Add `km_cut` argument to analysis.
 - converter: Add operator checking before optimization.
 - E2E simulator: Change dynasty library fetching method.
 - Minor bug fixes.
- **[v0.18.0]**
 - ONNX is updated to 1.7.0.
 - Introduce WebGUI.
 - Adjust 720 and 530 IP Evaluator default hardware specification.
 - Add more analysis options.
- For history versions, please check this link (<https://doc.kneron.com/docs/#toolchain/history>).

1. Installation

Review the system requirements below before start installing and using the toolchain.

1.1 System requirements

1. **Hardware:** Minimum quad-core CPU, 4GB RAM and 6GB free disk space.
2. **Operating system:** Window 10 x64 version 1903 or higher with build 18362 or higher. Ubuntu 16.04 x64 or higher. Other OS which can run docker later than 19.03 may also work. But they are not tested. Please take the risk yourself.
3. **Docker:** Docker Desktop later than 19.03. Here is a link (<https://www.docker.com/products/docker-desktop>) to download Docker Desktop.

TIPS:

For Windows 10 users, we recommend using docker with wsl2, which is Windows subsystem Linux provided by Microsoft. Here is how to install wsl2 (<https://docs.microsoft.com/en-us/windows/wsl/install-win10>) and how to install and run docker with wsl2 (<https://docs.docker.com/docker-for-windows/wsl/>). Also, you might want to adjust the resources docker use to ensure the tools' normal usage. Please check the FAQ at the end of this document on how to do that.

Please double-check whether the docker is successfully installed and callable from the console before going on to the next section. If there is any problem about the docker installation, please search online or go to the docker community for further support. The questions about the docker is beyond the reach of this document.

1.2 Pull the latest toolchain image

All the following steps are on the command line. Please make sure you have the access to it.

TIPS:

You may need `sudo` to run the docker commands, which depends on your system configuration.

You can use the following command to pull the latest toolchain docker.

```
docker pull kneron/toolchain:latest
```

Note that this document is compatible with toolchain v0.18.2. You can find the version of the toolchain in `/workspace/version.txt` inside the docker. If you find your toolchain is later than v0.18.2, you may need to find the latest document from the online document center (<http://doc.kneron.com/docs>).

2. Toolchain Docker Overview

After pulling the desired toolchain, now we can start walking through the process. In all the following sections, we use `kneron/toolchain:latest` as the docker image. Before we actually start the docker, we'd better provide a folder which contains the model files you want to test in our docker, for example, `/mnt/docker`. Then, we can use the following command to start the docker and work in the docker environment:

```
docker run --rm -it -v /mnt/docker:/docker_mount kneron/toolchain:latest
```

TIPS:

The mount folder path here is recommended to be an absolute path.

Here are the brief explanations for the flags. For detailed explanations, please visit docker documents (<https://docs.docker.com/engine/reference/run/>).

- `--rm`: the container will be removed after it exists. Each time we use `docker run`, we create a new docker container. Thus, without this flag, the docker will consume more and more disk space.
- `-it`: enter the interactive mode so we can use the bash.
- `-v`: mount a folder into the docker container. Thus, we can visit the desired files from the host and save the result from the container.

2.1 Folder structure

After logging into the container, you are under `/workspace`, where all the tools are. Here is the folder structure and their usage:

```
/workspace
|-- E2E_Simulator      # End to end simulator
|-- ai_training        # AI training project.
|-- cmake              # Environment
|-- examples           # Example for the workflow, will be used later.
|-- libs               # The libraries
|  |-- ONNX_Convertor  # ONNX Converters and optimizer scripts, will be discussed :
|  |-- compiler        # Compiler for the hardware and the IP evaluator to infer th
|  |-- dynasty         # Simulator which only simulates the calculation.
|  |-- fpAnalyser      # Analyze the model and provide fixed point information.
|  |-- hw_c_sim        # Hardware simulator which simulate all the hardware behavio
|-- miniconda          # Environment
|-- scripts            # Scripts to run the tools, will be discussed in section 3.
`-- version.txt
```

2.2 Work flow

Before we start actually introducing the usage, let us go through the general work flow.

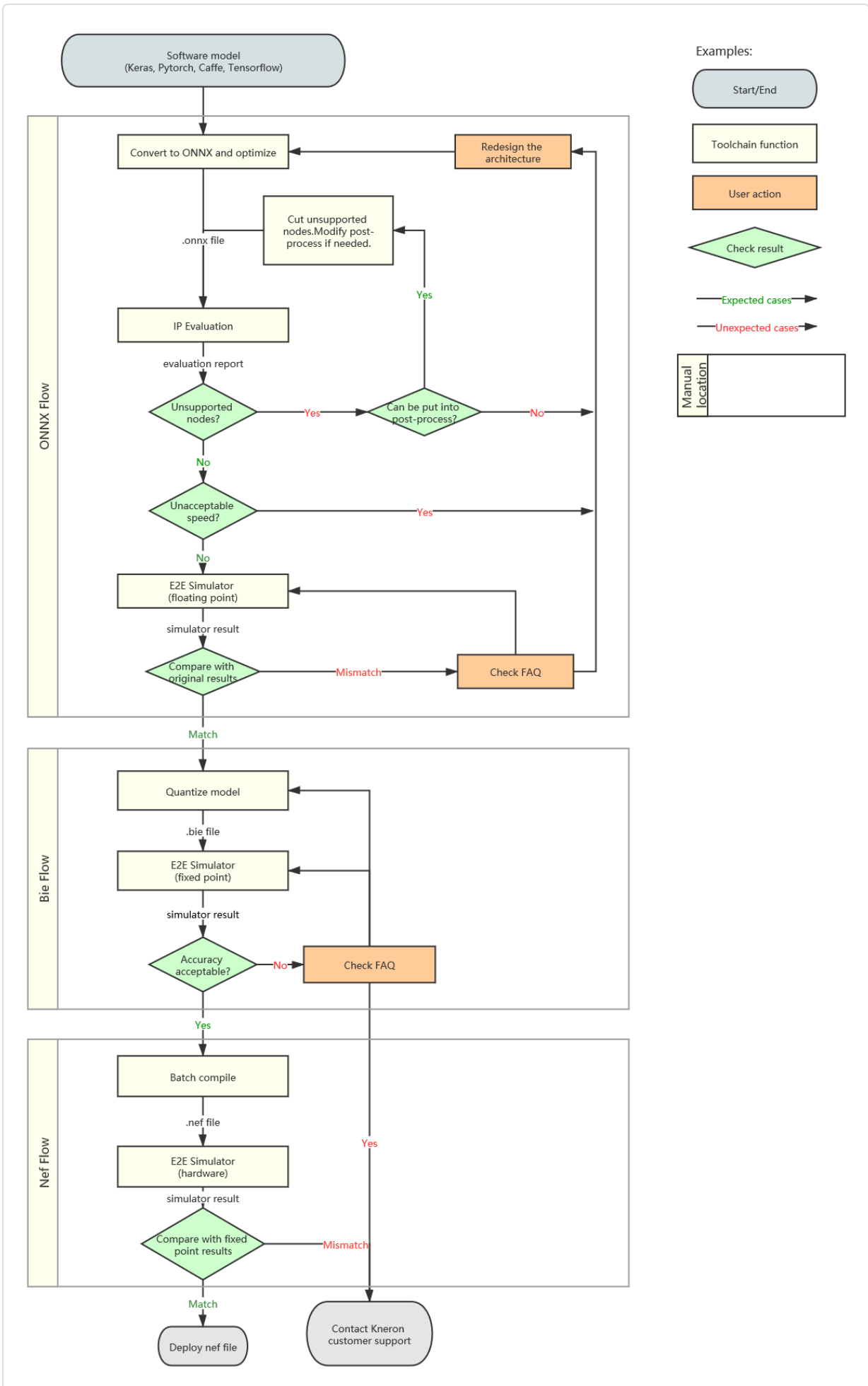


Figure 1. Diagram of working flow

To keep the diagram as clear as possible, some details are omitted. But it is enough to show the general steps. There are three main sections:

1. ONNX section. Convert the model from different platforms to onnx and optimize the onnx file. Evaluate the onnx the model to check the operator support and the estimate performance. Then, test the onnx model and compare the result with the source.
2. Bie section. Quantize the model and generate bie file. Test the bie file and compare the result with the previous step.
3. Nef section. Batch compile multiple bie models into a bie binary file. Test the nef file and compare the result with the previous step.

The workflow has been changed a lot since v0.15.0. We recommend Python API instead of the original scrips for a more smooth workflow. But this doesn't means the scrips used before are abandoned. They are still available and the document can be found in Command Line Script Tools (http://doc.kneron.com/docs/toolchain/command_line/). The detailed document of the Python API can be found in Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/)

2.3 Supported operators

Table 1.1 shows the list of functions KDP520 supports base on ONNX 1.6.1.

Table 1.1 The functions KDP520 NPU supports

Type	Operarots	Applicable Subset	Spec.
Convolution	Conv	Kernel dimension Strides	1x1 up to 11x11 1,2,4
	Pad		0-15
	Depthwise Conv		Yes
	Deconvolution		Use Upsampling + Conv
Pooling	MaxPool	3x3	stride 1,2,3
	MaxPool	2x2	stride 1,2
	AveragePool	3x3	stride 1,2,3
	AveragePool	2x2	stride 1,2
	GlobalAveragePool		support
Activation	GlobalMaxPool		support
	Relu		support
	LeakyRelu		support
Other processing	PRelu		support
	BatchNormalization		support

Type	Operarots	Applicable Subset	Spec.
	Add		support
	Concat		axis = 1
	Gemm or Dense/Fully Connected		support
	Flatten		support
	Clip		min = 0

Table 1.2 shows the list of functions KDP720 supports base on ONNX 1.6.1.

Table 1.2 The functions KDP720 NPU supports

Node	Applicable Subset	Spec.
Relu		support
PRelu		support
LeakyRelu		support
Sigmoid		support
Clip		min = 0
Tanh		support
BatchNormalization	up to 4D input	support
Conv		strides < [4, 16]
Pad		spacial dimension only
ConvTranspose		strides = [1, 1], [2, 2]
Upsample		support
Gemm	2D input	support
Flatten	Before Gemm	support
Add		support
Concat		axis = 1
Mul		support
MaxPool		kernel = [1, 1], [2, 2], [3, 3]
AveragePool	3x3	kernel = [1, 1], [2, 2], [3, 3]
GlobalAveragePool	4D input	support
GlobalMaxPool		support
MaxRoiPool		support
Slice	input dimension <= 4	support

3 ONNX Workflow

Our toolchain utilities take ONNX files as inputs. The ONNX workflow is mainly about convert models from other platforms to ONNX and prepare the onnx for the quantization and compilation. There are three main steps: conversion, evaluation and testing.

3.0 Preparation

Most of the following codes are using the python API. You may need to import necessary libraries before copying the code in this document.


```
import onnx
import ktc
```

The details of the Python API can be found in Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/) document.

3.1 Model Conversion and Optimization

The onnx converter part currently support Keras, TFLite, a subset of Tensorflow, Caffe and Pytorch. Here, we will only briefly introduce some common usage of the python API. This part of python API is based on our converter and optimizer open-source project which can be found on Github https://github.com/kneron/ONNX_Convertor (https://github.com/kneron/ONNX_Convertor). The detailed usage of those converter scripts can be found in ONNX Converter (<http://doc.kneron.com/docs/toolchain/converters/>). The Tensorflow ktc api is not introduced here. We recommend export the tensorflow model to tflite and convert the tflite model. If you really want to try convert a pb file, please check the onnx converter project.

The example models used in the following converter command are not included in the docker by default. They can be found on Github <https://github.com/kneron/ConvertorExamples> (<https://github.com/kneron/ConvertorExamples>). You can download them through these terminal commands:

```
git clone https://github.com/kneron/ConvertorExamples.git
cd ConvertorExamples && git lfs pull
```

3.1.1 Keras to ONNX

For Keras, our converter support models from Keras 2.2.4. **Note that `tf.keras` and Keras 2.3 is not supported.** You may need to export the model as tflite model and see section 3.1.4 for TF Lite model conversion.

Suppose there is an onet hdf5 model exported By Keras, you need to convert it to onnx by the following python code:

```
result_m = ktc.onnx_optimizer.keras2onnx_flow('/data1/ConvertorExamples/keras_examp
```

In this line of python code, `ktc.onnx_optimizer.keras2onnx_flow` is the function that takes an hdf5 file path and convert the hdf5 into an onnx object. The return value `result_m` is the converted onnx object. It need to take one more optimization step (section 3.1.5) before going into the next section.

There might be some warning log printed out by the Tensorflow backend, but we can ignore it since we do not actually run it. You can check whether the conversion succeed by whether there is any exception raised.

This function has more parameters for fine-tuning. Please check Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/) if needed.

3.1.2 Pytorch to ONNX

Our Python API do not actually convert the python model. It only takes Pytorch exported onnx object as the input and optimize it. Please checkout the tips below on how to export an onnx with Pytorch.

The Pytorch inside the docker is version 1.7.1. We currently only support models exported by Pytorch version $\geq 1.0.0$, $\leq 1.7.1$. Other versions are not tested.

TIPS on export onnx using `torch.onnx` :

You can use `torch.onnx` to export your model into onnx object. Here is the Pytorch to ONNX document (<https://pytorch.org/docs/stable/onnx.html#example-alexnet-from-pytorch-to-onnx>) for `onnx.torch`. An example code for exporting the model is:

```
import torch.onnx
dummy_input = torch.randn(1, 3, 224, 224)
torch.onnx.export(model, dummy_input, 'output.onnx', opset_version=11)
```

In the example, `(1, 3, 224, 224)` are batch size, input channel, input height and input width. `model` is the pytorch model object you want to export. `output.onnx` is the output onnx file path.

Pytorch exported onnx needs to pass through a special optimization designed for pytorch exported models first. Suppose the input file is loaded into a onnx object `exported_m`, here is the python code for pytorch exported onnx optimization:

```
result_m = ktc.onnx_optimizer.torch_exported_onnx_flow(exported_m)
```

In this line of python code, `ktc.onnx_optimizer.torch_exported_onnx_flow` is the function that takes an onnx object and optimize it. The return value `result_m` is the optimized onnx object. It need to take one more general onnx optimization step (section 3.1.5) before going into the next section.

For the example in ConverterExamples project, the whole process would be:

```

import torch
import torch.onnx

# Load the pth saved model
pth_model = torch.load("/data1/ConvertorExamples/keras_example/resnet34.pth", map_lo
# Export the model
dummy_input = torch.randn(1, 3, 224, 224)
torch.onnx.export(pth_model, dummy_input, '/data1/resnet34.onnx', opset_version=11)
# Load the exported onnx model as an onnx object
exported_m = onnx.load('/data1/resnet34.onnx')
# Optimize the exported onnx object
result_m = ktc.onnx_optimizer.torch_exported_onnx_flow(exported_m)

```

There might be some warning log printed out by the Pytorch because our example model is a little old. But we can ignore it since we do not actually run it. You can check whether the conversion succeed by whether there is any exception raised.

Crash due to name conflict

If you meet the errors related to `node not found` or `invalid input`, this might be caused by a bug in the onnx library. Please try set `disable_fuse_bn` to `True`. The code would be:

```

result_m = ktc.onnx_optimizer.torch_exported_onnx_flow(exported_m, disable_fuse_bn=

```

3.1.3 Caffe to ONNX

For caffe, we only support model which can be loaded by Intel Caffe 1.0 (<https://github.com/intel/caffe>).

Here we will use the example from ConvertorExamples. You can find two files for the caffe model: the model structure definition file `mobilenetv2.prototxt` and the model weight file `mobilenetv2.caffemodel`. Here is the example python code for model conversion:

```

result_m = ktc.onnx_optimizer.caffe2onnx_flow('/data1/ConvertorExamples/caffe_examp

```

In this line of python code, `ktc.onnx_optimizer.caffe2onnx_flow` is the function that takes caffe model file paths and convert the them into an onnx object. The return value `result_m` is the converted onnx object. It need to take one more optimization step (section 3.1.5) before going into the next section.

3.1.4 TF Lite to ONNX

We only support unquantized TF Lite models for now. Also tensorflow 2 is not supported yet.

Suppose we are using the tflite file `model_unquant.tflite` from the ConvertorExamples, here is the example python code:

```
result_m = ktc.onnx_optimizer.tflite2onnx_flow('/data1/ConvertorExamples/tflite_exar
```

In this line of python code, `ktc.onnx_optimizer.tflite2onnx_flow` is the function that takes an tflite file path and convert the tflite into an onnx object. The return value `result_m` is the converted onnx object. It need to take one more optimization step (section 3.1.5) before going into the next section.

There might be some warning log printed out by the Tensorflow backend, but we can ignore it since we do not actually run it. You can check whether the conversion succeed by whether there is any exception raised.

This function has more parameters for fine-tuning. Please check Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/) if needed.

3.1.5 ONNX Optimization

We provide a general onnx optimize API. We strongly recommend that all the onnx, including the onnx generated from the previous subsections, shall pass this API before going into the next section. This general onnx optimization API would modify the onnx graph to fit the toolchain and Kneron hardware specification. The optimization includes: inference internal value_info shapes, fuse consecutive operators, eliminate do-nothing operators, replace high-cost operators with low-cost operators, etc..

Suppose we have a onnx object, here is the example python code:

```
optimized_m = ktc.onnx_optimizer.onnx2onnx_flow(result_m, eliminate_tail=True, opt_r
```

In this line of python code, `ktc.onnx_optimizer.onnx2onnx_flow` is the function that takes an onnx object and optimize it. The return value `result_m` is the converted onnx object.

Note that for hardware usage, `eliminate_tail` should be set to true as in the example. This option eliminate the no calculation operators and the npu unsupported operators (Reshape, Transpose, ...) at the last of the graph. However, since this changes the graph structure, you may need to check the model yourself and add the related functions into post-process to keep the algorithm consistent. If you only want to use onnx model for software testing, the `eliminate_tail` can be set to false to keep the model same as the input from the mathematics perspective. `opt_matmul` is also for hardware usage which optimize the MatMul nodes according to the NPU limit. By default, it is set to False. You only need to enable this flag if there are MatMul nodes with inputs more than 2D.

This function has more parameters for fine-tuning. Please check Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/) if needed.

By the way, to save the model, you can use the following function from the onnx package.

```
# Save the onnx object optimized_m to path /data1/optimized.onnx.  
onnx.save(optimized_m, '/data1/optimized.onnx')
```

Crash due to name conflict

If you meet the errors related to `node not found` or `invalid input`, this might be caused by a bug in the onnx library. Please try set `disable_fuse_bn` to `True`. The code would be:

```
optimized_m = ktc.onnx_optimizer.onnx2onnx_flow(result_m, eliminate_tail=True, disal
```

Error due to the opset version

If you have met errors which are related to the opset version or the ir version. Please check section 3.1.6 to update your model first.

3.1.6 ONNX Opset Upgrade

From toolchain version 0.14.0, ONNX in the docker has been updated from 1.4.1 to 1.6. And the default opset that converters support is changed from opset 9 into opset 11. IR version is updated from 4 to 6. Thus, if you have a onnx model with opset 9 or IR version 4, you may need to update it with the following Python API:

```
new_m = ktc.onnx_optimizer.onnx1_4to1_6(old_m)
```

In this line of python code, `ktc.onnx_optimizer.onnx1_4to1_6` is the function that takes an old version onnx object and upgrade it. The return value `new_m` is the converted onnx object. It need to take one more optimization step (section 3.1.5) before going into the next section. Even if you have already passed optimizar before, we still recommend you do it again after this upgrade.

3.1.7 ONNX Editor

KL520/KL720/KL530 NPU supports most of the compute extensive OPs, such as Conv, BatchNormalization, Fully Connect/GEMM, in order to speed up the model inference run time. On the other hand, there are some OPs that KL520 NPU cannot support well, such as `Softmax` or `Sigmoid`. However, these OPs usually are not compute extensive and they are better to execute in CPU. Therefore, Kneron provides python APIs which help user modify the model so that KL520 NPU can run the model more efficiently.

You can find the detailed description of this tool from Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/) for the python API and ONNX Converter (<http://doc.kneron.com/docs/toolchain/converters/>) for the command usage.

3.2 IP Evaluation

Before we start quantizing the model and try simulating the model, we need to test if the model can be taken by the toolchain structure and estimate the performance. IP evaluator is such a tool which can estimate the performance of your model and check if there is any operator or structure not supported by our toolchain.

From this section, we would use the LittleNet which is included in the docker as an example. You can find it under `/workspace/examples/LittleNet`.

We need to create a `ktc.ModelConfig` object. The `ktc.ModelConfig` is the class which contains the basic needed information of a model. You can initialize it through the API below.

```
km = ktc.ModelConfig(id, version, platform, onnx_model=None, onnx_path=None, bie_pa
```

- `id` is the identifier of the model. It should be a integer greater than 0. ID before 32768 are reserved for Kneron models. Please use ID greater than 32768 for custom models.
- `version` is the model version. It should be a four digit hex code which is written as string, e.g. '001a'.
- `platform` is the target platform for this model. It should be '520', '720' or '530'.
- `onnx_model`, `onnx_path` and `bie_path`. User should provide one of those three parameter and only one. It stores the model itself. `onnx_model` takes the onnx object which is generated through the converter APIs or loaded through onnx library. `onnx_path` is the path to a onnx file. `bie_path` is the path to a bie file. The bie file is the file generated by the kneron toolchain after quantization, which is introduced in the later section.

For this example, we create the LittleNet ModelConfig with the following python code:

```
km = ktc.ModelConfig(32769, "0001", "520", onnx_path="/workspace/examples/LittleNet,
```

`evaluate` is class function of the `ktc.ModelConfig`.

```
eval_result = km.evaluate()
```

The evaluation result will be returned as string. User can also find the evaluation result under `/data1/compiler/`. But the report file names are different for different platforms.

- 520: `ip_eval_prof.txt`
- 720, 530: `ProfileResult.txt`

If the model is not supported, there would be warning messages or exceptions. Please modify the model structure referring to the message. Please check the report to see if the performance meets your expectation. Please consider redesign the network structure. Also note that the evaluator report only considers the performance of NPU. Thus, if the model contains many operators that are not supported by NPU but by CPU, the actual performance would be even worse.

TIPS: You can find the profiling configuration under `/workspace/scripts/res`. The configuration files are named like `ip_config_<platform>.json`. You can change the bandwidth according to your scenario .

3.3 E2E Simulator Check (Floating Point)

Before going into the next section of quantization, we need to ensure the optimized onnx file can produce the same result as the originally designed model.

Here we introduce the E2E simulator which is the abbreviation for end to end simulator. It can inference a model and simulate the calculation of the hardware. The inference function of the E2E simulator is called `ktc.kneron_inference`. Here we are using the onnx as the input model. But it also can take bie file and nef file which would be introduced later.

The python code would be like:

```
inf_results = ktc.kneron_inference(input_data, onnx_file="/workspace/examples/LittleNet.onnx")
```

In the code above, `inf_results` is a list of result data. `onnx_file` is the path to the input onnx. `input_data` is a list of input data after preprocess the `input_names` is a list of string mapping the input data to specific input on the graph using the sequence.

Note that the input should have the same dimension as the model but in channel last format.

Here we provide a very simple preprocess function which only do the resize and normalization:

```
from PIL import Image
import numpy as np

def preprocess(input_file):
    image = Image.open(input_file)
    image = image.convert("RGB")
    img_data = np.array(image.resize((112, 96), Image.BILINEAR)) / 255
    img_data = np.transpose(img_data, (1, 0, 2))
    return img_data
```

Then, the full code of this section would be:

```
from PIL import Image
import numpy as np

def preprocess(input_file):
    image = Image.open(input_file)
    image = image.convert("RGB")
    img_data = np.array(image.resize((112, 96), Image.BILINEAR)) / 255
    img_data = np.transpose(img_data, (1, 0, 2))
    return img_data

input_data = [preprocess("/workspace/examples/LittleNet/pytorch_imgs/Abdullah_0001.jpg")]
inf_results = ktc.kneron_inference(input_data, onnx_file="/workspace/examples/LittleNet.onnx")
```

Since we want to focus on the toolchain usage here, we do not provide any postprocess. In reality, you may want to have your own postprocess function in Python, too.

After getting the `inf_results` and post-process it, you may want to compare the result with the one generated by the source model. For example, if the source model is from pytorch, you may want to try inference the source model using Pytorch with the same input image to see if the results match. If the result mismatch, please check FAQ 1 for possible solution.

Since we do not actually has any source model here for the simplicity of example, we would skip the step of comparing result.

4 BIE Workflow

As mentioned briefly in the previous section, the bie file is the model file which is usually generated after quantization. It is encrypted and not available for visuanlization. In this chapter, we would go through the steps of quantization.

4.1 Quantization

Quantization is the step where the floating-point weight are quantized into fixed-point to reduce the size and the calculation complexity. The Python API for this step is called `analysis`. It is also a class function of `ktc.ModelConfig`. It takes a dictionary as input.

```
analysis(input_mapping, output_bie = None, threads = 4, mode=1)
```

- `input_mapping` is the a dictionary which maps a list of input data to a specific input name. Generally speaking, the quantization would be preciser with more input data.
- `output_bie` is the path where you want your bie generated. By default, it is under `/data1/fpAnalyser`.
- `threads` is the threads number you want to utilize. Please note more threads would lead to more RAM usage as well.
- `mode` is an optional flag to determine whether to skip the model verification step while doing the quantization. The model verification makes sure your model can be processed correctly by our toolchain. But this step could take more time and consume more system resources. Note that if your memory is not enough, the utility would raise segmentation fault. By default, this flag is set to 1, which means the verification is skipped.
 - 1: only analysis. Skip verification.
 - 2: Verification with one image.
 - 3: Verification with all provided images. (WARNING: This option takes very long time.)
- The return value is the generated bie path.

This is a very simple example usage. There are many more parameters for fine-tuning. Please check Please check Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/) if needed.

Please also note that this step would be very time-consuming since it analysis the model with every input data you provide.

Here as a simple example, we only use four input image as exmaple and run it with the `ktc.ModelConfig` object `km` created in section 3.2:

```
# Preprocess images and create the input mapping
input_images = [
    preprocess("/workspace/examples/LittleNet/pytorch_imgs/Abdullah_0001.png"),
    preprocess("/workspace/examples/LittleNet/pytorch_imgs/Abdullah_0002.png"),
    preprocess("/workspace/examples/LittleNet/pytorch_imgs/Abdullah_0003.png"),
    preprocess("/workspace/examples/LittleNet/pytorch_imgs/Abdullah_0004.png"),
]
input_mapping = {"data_out": input_images}

# Quantization
bie_path = km.analysis(input_mapping, output_bie = None, threads = 4)
```

This function has more parameters for fine-tuning. Please check Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/) if needed.

4.2 E2E Simulator Check (Fixed Point)

Before going into the next section of compilation, we need to ensure the quantized model do not lose too much precision.

We would use `ktc.kneron_inference` here, too. But here we are using the generated bie file as the input.

The python code would be like:

```
fixed_results = ktc.kneron_inference(input_data, bie_file=bie_path, input_names=["d
```

The usage is almost the same as using onnx. In the code above, `inf_results` is a list of result data. `bie_file` is the path to the input bie. `input_data` is a list of input data after preprocess the `input_names` is a list of model input name. The requirement is the same as in section 3.3. If your platform is not 520, you may need an extra parameter `platform`, e.g. `platform=720` or `platform=530`.

As mentioned above, we do not provide any postprocess. In reality, you may want to have your own postprocess function in Python, too.

After getting the `fixed_results` and post-process it, you may want to compare the result with the `inf_results` which is generated in section 3.3 to see if the precision lose too much. If the result is unacceptable, please check FAQ 2 for possible solutions.

5 NEF Workflow

The nef file is the binary file after compiling and can be taken by the Kneron hardware. But unlike the utilities mentioned above, the process in this chapter can compile multiple models into one nef file. This process is called batch process. In this chapter, we would go through how to batch compile and how to verify the nef file.

5.1 Batch Compile

Batch compile turns multiple models into a single binary file. But, we would start with single model first.

The Python API is very simple:

```
compile_result = ktc.compile(model_list)
```

The `compile_result` is the path for the generated nef file. By default, it is under `/data1/batch_compile`. It takes a list of `ktc.ModelConfig` object as the input `model_list`. The usage of `ktc.ModelConfig` can be found in section 3.2. Note that the `ModelConfig` object must have bie file inside. In details, it must be under either of the following status: the `ModelConfig` is initialized with `bie_path`, the `ModelConfig` is initialized with `onnx_model` or `onnx_path` but it has successfully run `analysis` function.

For the LittleNet example, please check the code below. Note that `km` is the `ktc.ModelConfig` object we generate in section 3.2 and use in the section 4.

```
compile_result = ktc.compile([km])
```

For multiple models, we can simply extend the model list.

```
# dummy.bie is not a real example bie which is available in the docker. Just for con
# Please adjust the parameters according to your actual input.
km2 = ktc.ModelConfig(32770, "0001", "520", bie_path="dummy.bie")
compile_result = ktc.compile([km, km2])
```

Note that for multiple models, all the models should share the same hardware platform and the model ID should be different.

This function has more parameters for fine-tuning. Please check Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/) if needed.

5.2 E2E Simulator Check (Hardware)

After compilation, we need to check if the nef can work as expected.

We would use `ktc.kneron_inference` here again. And we are using the generated nef file this time.

For the batch compile with only LittleNet, the python code would be like:

```
hw_results = ktc.kneron_inference(input_data, nef_file=compile_result)
```

The usage is a little different here. In the code above, `hw_results` is a list of result data. `nef_file` is the path to the input nef. `input_data` is a list of input data after preprocess. The requirement is the same as in section 3.3. If your platform is not 520, you may need an extra parameter `platform`, e.g. `platform=720` or `platform=530`.

Here we use the same input `input_data` which we used in section 3.3. And the `compile_result` is the one that generated with only LittleNet model.

As mentioned above, we do not provide any postprocess. In reality, you may want to have your own postprocess function in Python, too.

For nef file with multiple models, we can specify the model with model ID:

```
hw_results = ktc.kneron_inference(input_data, nef_file=compile_result, model_id=32769)
```

After getting the `hw_results` and post-process it, you may want to compare the result with the `fixed_results` which is generated in section 4.2 to see if the results match. If the results mismatch, please contact us directly through forum <https://www.kneron.com/forum/> (<https://www.kneron.com/forum/>).

5.3 NEF Combine

This section is not part of the normal workflow. But it would be very useful when you already have multiple nef files, some of which might from different versions of the toolchain, and you want to combine them into one. Here we provide a python API to achieve it.

```
ktc.combine_nef(nef_path_list, output_path = "/data1/combined")
```

Here the `nef_path_list` shall be a list of the `str` which are the path to the nef files you want to combine. It should not be empty. And the second argument is optional. It should be the output folder path of the combined nef. By default, it should be `/data1/combined`. The return value is the output folder path. The combined nef file would be under the output folder and be named as `models_<target>.def`. For example, if your target platform is 520, the result file name would be `models_520.nef` inside the output folder.

Here is an usage example. Suppose you have three 520 nef files: `/data1/model_32769.nef`, `/data1/model_32770.nef` and `/data1/model_32771.nef`. Then you can use the following code to combine them. The result is `"/data1/combined/models_520.nef"`

```
ktc.combine_nef(['/data1/model_32769.nef', '/data1/model_32770.nef', '/data1/model_32771.nef'], output_path="/data1/combined")
```

6 What's Next

- Check `/workspace/examples/LittleNet/python_api_workflow.py`. All the example Python API usages are inside this runnable script.
- Try it with your own models.
- Check the YOLO Example (http://doc.kneron.com/docs/toolchain/yolo_example/) for a step-by-step walk through using YOLOv3 as the example.
- Check the Toolchain Python API (http://doc.kneron.com/docs/toolchain/python_api/) document for more detailed Python API usage.
- Check the Command Line Script Tools (http://doc.kneron.com/docs/toolchain/command_line/) for command line script usage.
- Check the ONNX Converter (<http://doc.kneron.com/docs/toolchain/converters/>) for the usage of underlying project https://github.com/kneron/ONNX_Convertor (https://github.com/kneron/ONNX_Convertor).
- Check the Web GUI (http://doc.kneron.com/docs/toolchain/toolchain_webgui/) for a simple web interface.

FAQ

1. What if the E2E simulator results from the original model and the optimized onnx mismatch?

Please double check if the final layers are cut due to unsupported by NPU. If so, please add the deleted operator as part of the E2E simulator post process and test again.

Otherwise, please search on forum <https://www.kneron.com/forum/categories/ai-model-migration> (<https://www.kneron.com/forum/categories/ai-model-migration>). You can also contact us through the forum if no match issue found. The technical support would reply directly to your post.

2. What if the E2E simulator results of floating-point and fixed-point lost too match accuracy?

Please try the following solutions:

1. Try redoing the analysis with more image that are the expected input of the network.
2. Double check if the cut final CPU nodes are added in post-process.
3. Fine tuning the analysis with outlier and quantize mode.

If none of the above works, please search on forum <https://www.kneron.com/forum/categories/ai-model-migration> (<https://www.kneron.com/forum/categories/ai-model-migration>). You can also contact us through the forum if no match issue found. The technical support would reply directly to your post.

3. How to adjust the system resources usage of the docker?

To ensure the quantization tool can work, we recommend the docker has at least 4GB of memory. The actual required size depends on your model size and the image number of quantization.

For Linux users, by default, docker can share all the CPU and memory resources of their host machine. So, this isn't a problem. But for Windows users, not like Linux, the system resources are not shared. User might want to adjust the resources usage by themselves.

For the docker based on wsl2, as we recommended in the section 1 of this document, it can use up to 50% of your total system memory and all the CPU resources. And here is an article introducing how to manage the system resources used by wsl2 (<https://ryanharrison.co.uk/2021/05/13/wsl2-better-managing-system-resources.html#:~:text=1%20Setting%20a%20WSL2%20Memory%20Limit.%20By%20default,the%20WSL2%20Virtual%20Disk.%20...%204%20Docker.%20>).

For the docker based on wsl, users can find the management of the system resources directly in the settings of the docker.

For the docker toolbox, it is actually based on the VirtualBox virtual machine. So, users need to find which virtual machine the docker is using first. Users need to start the docker terminal to ensure the docker is running before we start. And here is the following procedure:

- Open the VirtualBox management tool.

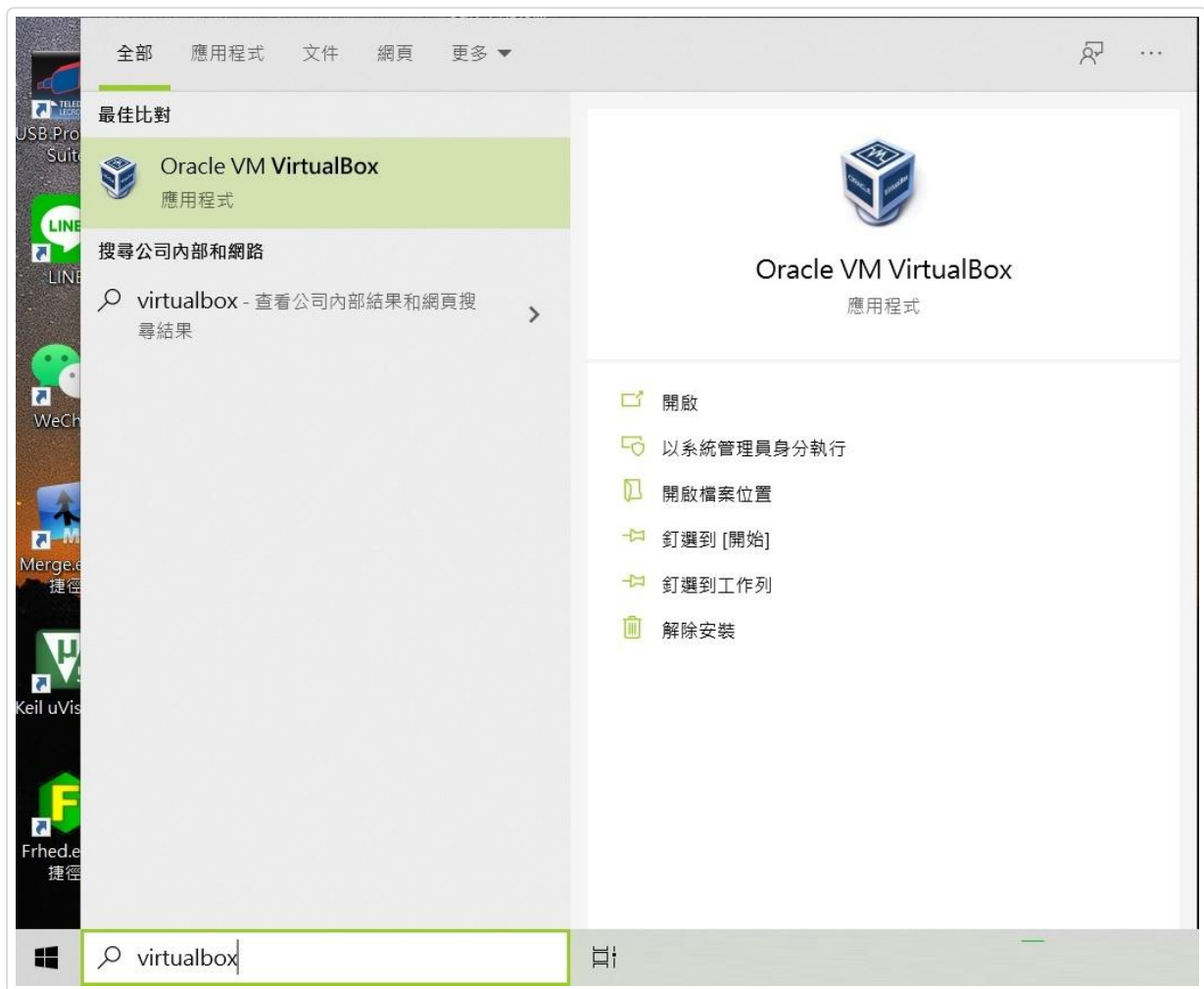


Figure FAQ3.1 VirtualBox

- Check the status. There should be only one virtual machine running if there is no other virtual machines started manually by the user.

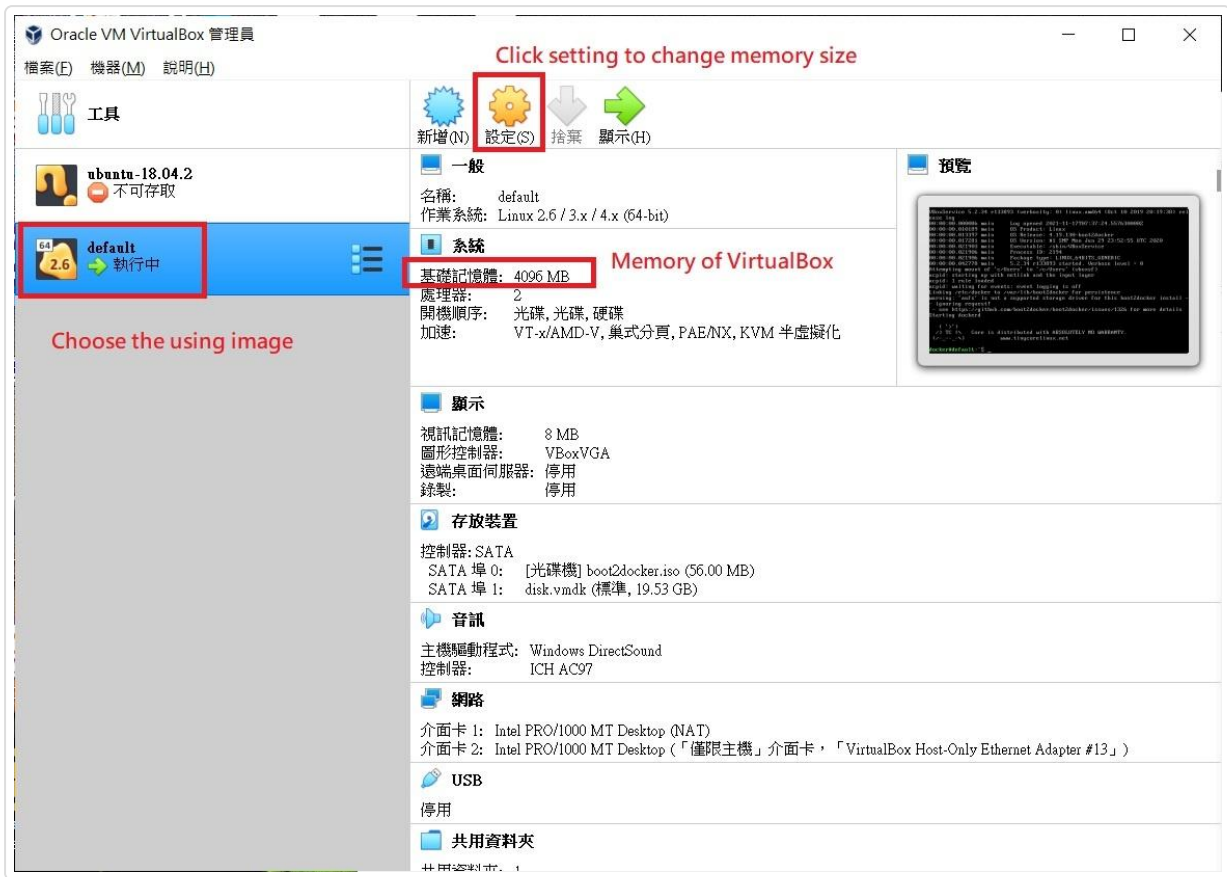


Figure FAQ3.2 VM status

- Close the docker terminal and shutdown the virtual machine before we adjust the resources usage.

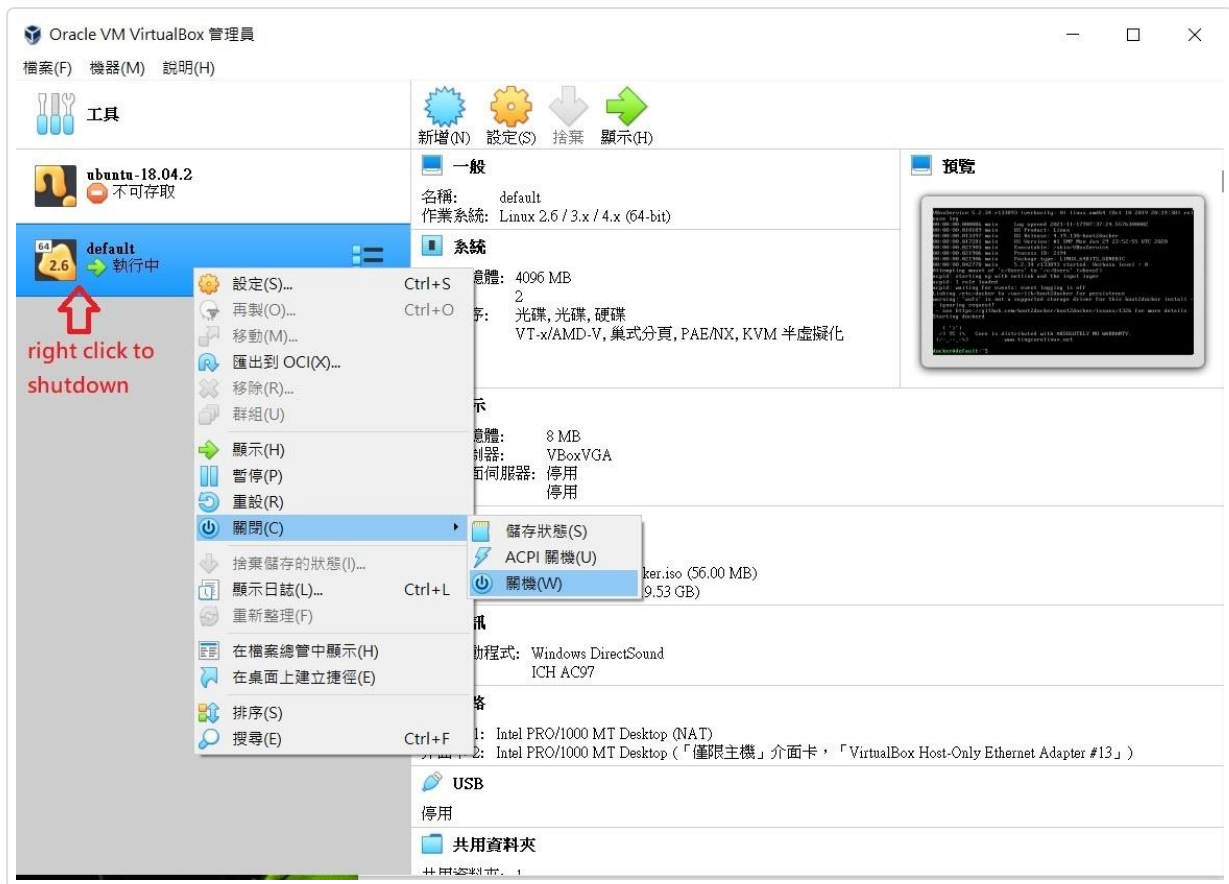


Figure FAQ3.3 VM shutdown

- Adjust the memory usage in the virtual machine settings. You can also change the cpu count here as well.

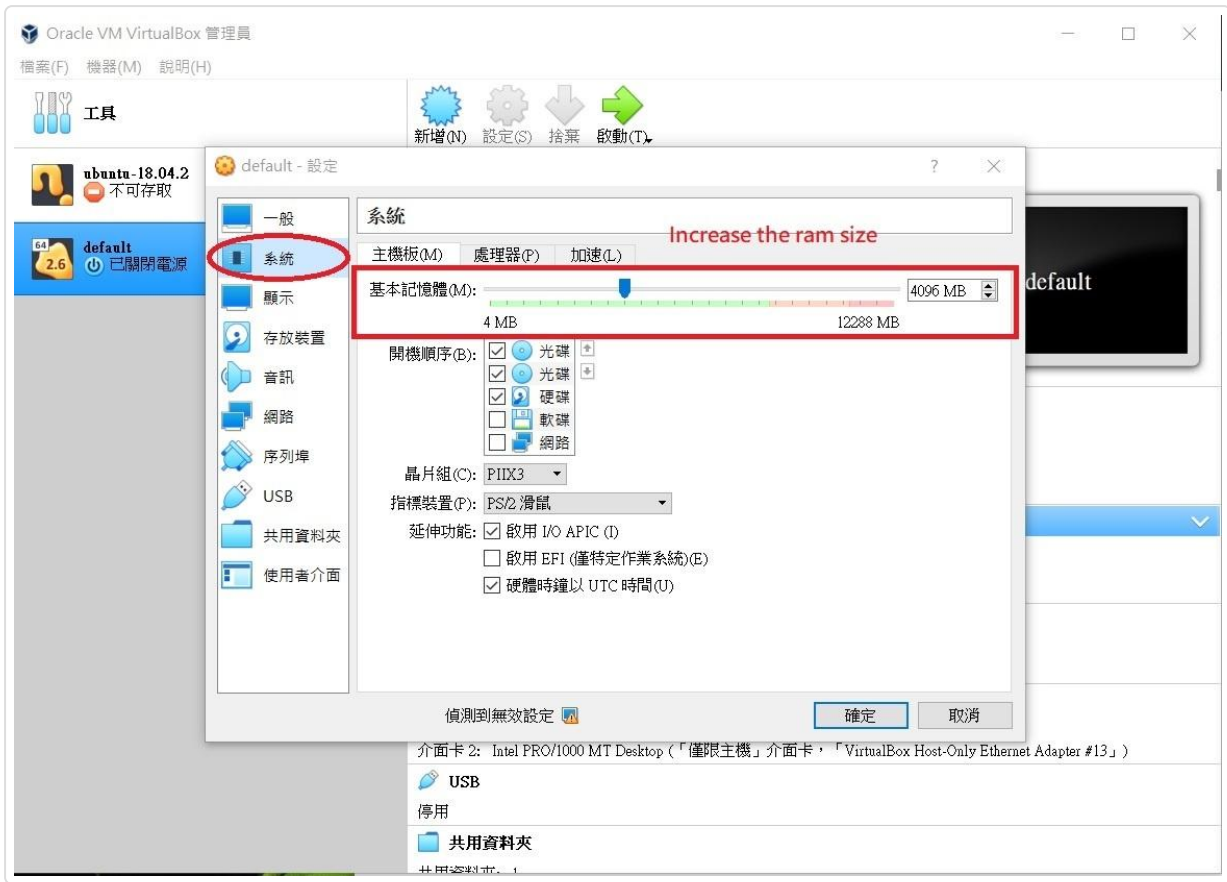


Figure FAQ3.4 VM settings

- Save the setting and restart the docker terminal. Now you can use more memory in your docker container.

`Systemctl refresh systemd_0.8.0-host_apt/systemd/`

`Yolo Example (../yolo_example/)`

[Edit on GitHub \(https://github.com/kneron/document_center/edit/master/docs/toolchain/manual.md\)](https://github.com/kneron/document_center/edit/master/docs/toolchain/manual.md)

Documentation built with MkDocs (<http://www.mkdocs.org/>) using Windmill (<https://github.com/gristlabs/mkdocs-windmill>) theme by Grist Labs.