




# ESP-DL User Guide




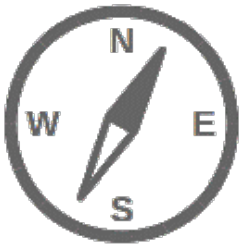

Release v3.0.0-101-gbf2ee2cbdc  
Espressif Systems  
Mar 12, 2025



# Table of contents

|  |          |
|--|----------|
| <b>Table of contents</b>                               | <b>i</b> |
| <b>1 Introduction</b>                                  | <b>3</b> |
| 1.1 Introduction                                       | 3        |
| 1.1.1 Overview   | 3        |
| 1.2 ESP-DL Project Organization                        | 4        |
| 1.2.1 <b>dl (Deep Learning)</b>                        | 4        |
| 1.2.2 <b>vision</b>                                    | 5        |
| 1.2.3 <b>fbs_loader (FlatBuffers Loader)</b>           | 5        |
| 1.2.4 <b>Other Files</b>                               | 5        |
| <b>2 Getting Started</b>                               | <b>7</b> |
| 2.1 Hardware Requirements                              | 7        |
| 2.2 Software Requirements                              | 7        |
| 2.2.1 ESP-IDF  | 7        |
| 2.2.2 ESP-PPQ  | 7        |
| 2.3 Quick Start  | 8        |
| 2.3.1 Example Compile & Flash                          | 8        |
| 2.3.2 Example Configuration                            | 8        |
| 2.3.3 Trouble shooting                                 | 8        |
| 2.4 Model Quantization                                 | 8        |
| 2.5 Model deployment                                   | 8        |
| <b>3 Tutorials</b>                                     | <b>9</b> |
| 3.1 How to quantize model                              | 9        |
| 3.1.1 Preparation                                      | 9        |
| 3.1.2 Pre-trained model                                | 9        |
| 3.1.3 Quantize and export .espdl                       | 10       |
| 3.2 How to load & test & profile model                 | 10       |
| 3.2.1 Preparation                                      | 11       |
| 3.2.2 Load model from rodata                           | 11       |
| 3.2.3 Load model from partition                        | 12       |
| 3.2.4 Load model from sdcard                           | 12       |
| 3.2.5 Test whether on-board model inference is correct | 13       |
| 3.2.6 Profile model memory usage                       | 13       |
| 3.2.7 Profile model inference latency                  | 13       |
| 3.3 How to run model                                   | 13       |
| 3.3.1 Preparation                                      | 14       |
| 3.3.2 Load model                                       | 14       |
| 3.3.3 Get model input/output.                          | 14       |
| 3.3.4 Quantize Input                                   | 14       |
| 3.3.5 Dequantize output                                | 15       |
| 3.3.6 Model Inference                                  | 15       |
| 3.4 How to deploy MobileNetV2                          | 16       |
| 3.4.1 Preparation                                      | 16       |
| 3.4.2 Model quantization                               | 16       |
| 3.4.3 Model deployment                                 | 24       |

|              |  |           |
|--------------|--|-----------|
| 3.5          | How to deploy YOLO11n . . . . .            | 24        |
| 3.5.1        | Preparation . . . . .                      | 25        |
| 3.5.2        | Model quantization . . . . .               | 25        |
| 3.5.3        | Model deployment . . . . .                 | 36        |
| 3.6          | Creating a New Module (Operator) . . . . . | 36        |
| 3.6.1        | Understand the Base Module Class . . . . . | 37        |
| 3.6.2        | Create a New Module Class . . . . .        | 37        |
| <b>4</b>     | <b>API Reference</b>                       | <b>39</b> |
| 4.1          | Tensor API Reference . . . . .             | 39        |
| 4.1.1        | Header File . . . . .                      | 39        |
| 4.1.2        | Classes . . . . .                          | 39        |
| 4.2          | Module API Reference . . . . .             | 44        |
| 4.2.1        | Header File . . . . .                      | 45        |
| 4.2.2        | Classes . . . . .                          | 45        |
| 4.2.3        | Header File . . . . .                      | 47        |
| 4.2.4        | Classes . . . . .                          | 48        |
| 4.3          | Model API Reference . . . . .              | 48        |
| 4.3.1        | Header File . . . . .                      | 49        |
| 4.3.2        | Macros . . . . .                           | 49        |
| 4.3.3        | Classes . . . . .                          | 49        |
| 4.3.4        | Header File . . . . .                      | 54        |
| 4.3.5        | Classes . . . . .                          | 54        |
| 4.3.6        | Header File . . . . .                      | 58        |
| 4.3.7        | Classes . . . . .                          | 58        |
| 4.4          | Fbs API Reference . . . . .                | 59        |
| 4.4.1        | Header File . . . . .                      | 59        |
| 4.4.2        | Classes . . . . .                          | 59        |
| 4.4.3        | Header File . . . . .                      | 60        |
| 4.4.4        | Classes . . . . .                          | 60        |
| <b>Index</b> |  | <b>67</b> |
| <b>Index</b> |  | <b>67</b> |

|   |   |   |
|---|---|---|
|  |  |  |
| <a href="#">Get Started</a>   | <a href="#">Tutorials</a>   | <a href="#">API Reference</a>   |



# Chapter 1

## Introduction

### 1.1 Introduction

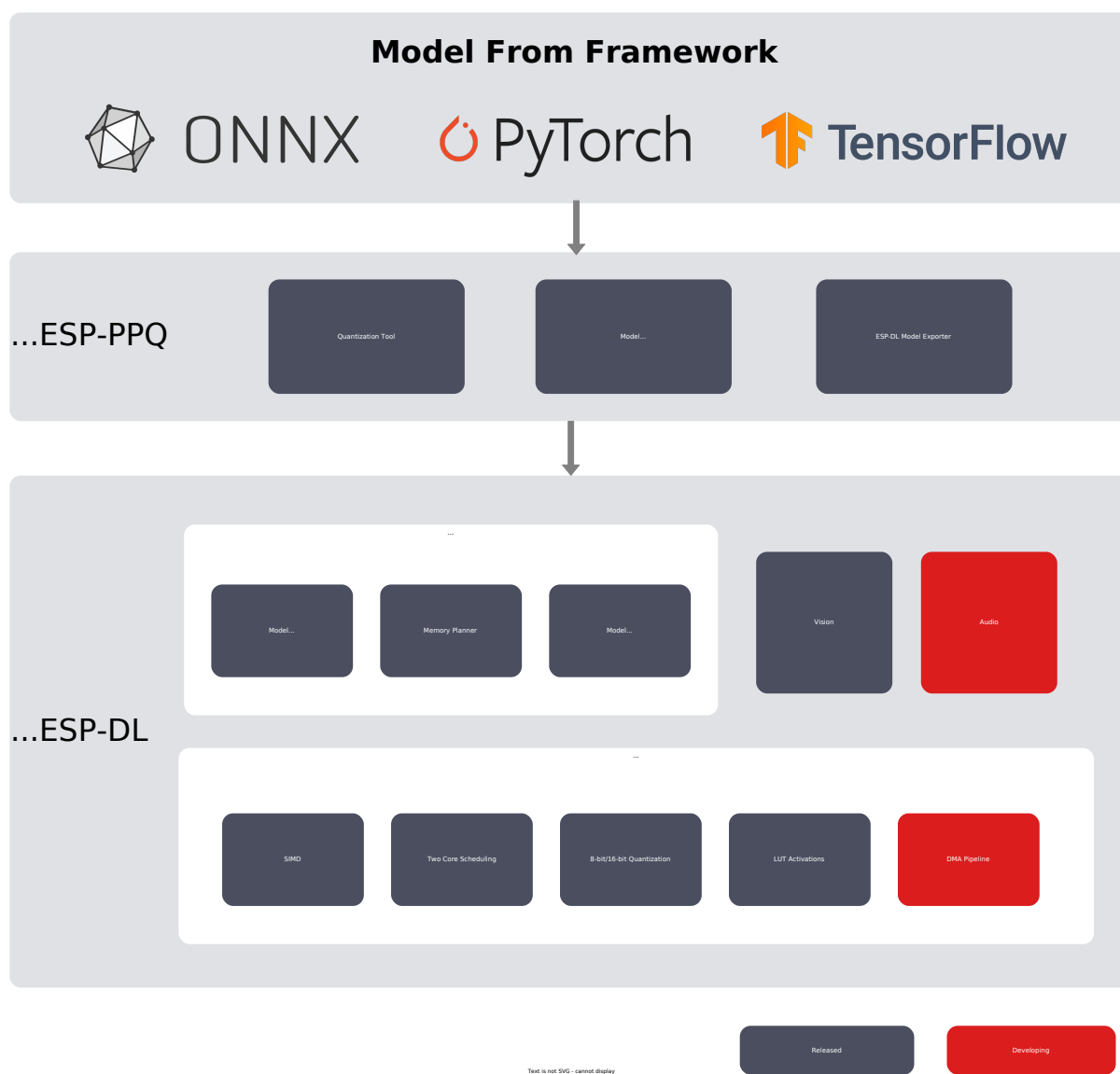
ESP-DL is a lightweight and efficient neural network inference framework designed specifically for ESP series chips. With ESP-DL, you can easily and quickly develop AI applications using Espressif's System on Chips (SoCs).

#### 1.1.1 Overview

ESP-DL offers APIs to load, debug, and run AI models. The framework is easy to use and can be seamlessly integrated with other Espressif SDKs. ESP-PPQ serves as the quantization tool for ESP-DL, capable of quantizing models from ONNX, Pytorch, and TensorFlow, and exporting them into the ESP-DL standard model format.

- **ESP-DL Standard Model Format:** This format is similar to ONNX but uses FlatBuffers instead of Protobuf, making it more lightweight and supporting zero-copy deserialization, with a file extension of *.espdl*.
- **Efficient Operator Implementation:** ESP-DL efficiently implements common AI operators such as Conv, Gemm, Add, and Mul. [The list of supported operators](#)
- **Static Memory Planner:** The memory planner automatically allocates different layers to the optimal memory location based on the user-specified internal RAM size, ensuring efficient overall running speed while minimizing memory usage.
- **Dual Core Scheduling:** Automatic dual-core scheduling allows computationally intensive operators to fully utilize the dual-core computing power. Currently, Conv2D and DepthwiseConv2D support dual-core scheduling.
- **8bit LUT Activation:** All activation functions except for ReLU and PReLU are implemented using an 8-bit LUT (Look Up Table) method in ESP-DL to accelerate inference. You can use any activation function, and their computational complexity remains the same.

The framework figures below illustrate the overall architecture of ESP-DL.



## 1.2 ESP-DL Project Organization

ESP-DL's modular design enables efficient development, maintenance, and scalability. The project is organized as follows:

### 1.2.1 dl (Deep Learning)

Core deep learning modules and tools, divided into submodules:

- **model** Loads, manages, and allocates memory for deep learning models. Includes `dl_model_base` and `dl_memory_manager`.
- **module** Interfaces for operations (convolution, pooling, activation, etc.). Files: `dl_module_conv`, `dl_module_pool`, `dl_module_relu` etc.
- **base** Implements operations for chips (esp32s3/esp32p4) with assembly support.
- **math** Mathematical operations (matrix functions). Files: `dl_math.hpp` and `dl_math.cpp`.
- **tool** Auxiliary functions (utility tools). Files: `dl_tool.hpp` and `dl_tool.cpp`.
- **tensor** Tensor classes. Files: `dl_tensor_base.hpp`.

## 1.2.2 vision

Computer vision modules divided into submodules:

- **classification** Image classification (model inference). Inference: `dl_cls_base`. Post-processors: `imagenet_cls_postprocessor`, `dl_cls_postprocessor`.
- **recognition** Feature extract (model inference). Feature database management (Enroll, delete, query). Pre-processor: `dl_feat_image_preprocessor`. Inference: `dl_feat_base`. Post-processor: `dl_feat_postprocessor`. Database: `dl_recognition_database`
- **image** Image process(resize, crop, warp affine). Color conversion(pixel, img). Image preprocessor (pipeline of resize, crop, color conversion, normalization, quantization). Image decoding/encoding (JPEG/BMP). Draw utility (point, hollow rectangle).  
Image process: `dl_image_process`. Color conversion: `dl_image_color`. Image preprocessor: `dl_image_preprocessor`. Image decoding/encoding: `dl_image_jpeg`, `dl_image_bmp`. Draw utility: `dl_image_draw`.
- **detect** Object detection (model inference). Inference: `dl_detect_base`. Post-processors: `dl_detect_msr_postprocessor`, `dl_detect_mnp_postprocessor`, `dl_detect_pico_postprocessor`.

## 1.2.3 fbs\_loader (FlatBuffers Loader)

Handles FlatBuffers models:

- **include** Headers: `fbs_loader.hpp`, `fbs_model.hpp`.
- **src** Implementations: `fbs_loader.cpp`.

## 1.2.4 Other Files

- **CMakeLists.txt** Project build configuration.
- **idf\_component.yml** Component metadata (name, version, dependencies).
- **README.md** Project documentation and usage.
- **LICENSE** License terms.





## Chapter 2

# Getting Started

### 2.1 Hardware Requirements

- An ESP32-S3 or ESP32-P4 development board. Recommended: ESP32-S3-EYE or ESP32-P4-Function-EV-Board
- PC (Linux)

---

**Note:** Some boards currently use Type C connectors. Make sure you use the right cable to connect the board!

---

### 2.2 Software Requirements

#### 2.2.1 ESP-IDF

ESP-DL runs based on ESP-IDF. For detailed instructions on how to get ESP-IDF, see the [ESP-IDF Programming Guide](#).

---

**Note:** Please use `release/v5.3` or higher version of [ESP-IDF](#).

---

#### 2.2.2 ESP-PPQ

ESP-PPQ is a quantization tool based on ppq. ESP-PPQ adds Espressif's customized quantizer and exporter based on [PPQ](#), which makes it convenient for users to select quantization rules that match ESP-DL according to different chip selections, and export them to standard model files that can be directly loaded by ESP-DL. ESP-PPQ is compatible with all PPQ APIs and quantization scripts. For more details, please refer to [PPQ documents and videos](#). If you want to quantize your own model, please install esp-ppq with the following command:

```
pip uninstall ppq
pip install git+https://github.com/espressif/esp-ppq.git
```

## 2.3 Quick Start

ESP-DL provides some out-of-the-box [examples](#)

### 2.3.1 Example Compile & Flash

```
idf.py set-target [Soc]
idf.py flash monitor -p [PORT]
```

Replace [Soc] with the specific chip, currently supports `esp32s3` and `esp32p4`.

### 2.3.2 Example Configuration

```
idf.py menuconfig
```

Some examples contain configurable options that can be configured using `idf.py menuconfig` after specifying the chip using `idf.py set-target`.

### 2.3.3 Trouble shooting

**Check ESP-IDF doc**

See [ESP-IDF DOC](#)

**Erase FLASH & Clear Example**

```
idf.py erase-flash -p [PORT]
```

Delete `build/`, `sdkconfig`, `dependencies.lock`, `managed_components/` and try again.

## 2.4 Model Quantization

First, please refer to [ESP-DL Operator Support State](#) to ensure that the operators in your model are supported.

ESP-DL must use the proprietary format `.espd1` for model deployment. Deep learning models need to be quantized and converted to the format before they can be used. ESP-PPQ provides two interfaces, `espd1_quantize_onnx` and `espd1_quantize_torch`, to support ONNX models and PyTorch models to be exported as `.espd1` models. Other deep learning frameworks, such as TensorFlow, PaddlePaddle, etc., need to convert the model to ONNX first. So make sure your model can be converted to ONNX model. For more details, please refer to:

- [How to quantize model](#)
- [How to quantize MobileNetV2](#)
- [How to quantize YOLO11n](#)

## 2.5 Model deployment

ESP-DL provides a series of APIs to quickly load and run models. For more details, see:

- [How to load & test & profile model](#)
- [How to run model](#)

# Chapter 3

## Tutorials

### 3.1 How to quantize model

ESP-DL must use a proprietary format `.espd1` for model deployment. This is a quantized model format that supports 8bit and 16bit. In this tutorial, we will take [quantize\\_sin\\_model](#) as an example to show how to use ESP-PPQ to quantize and export a `.espd1` model. The quantization method is Post Training Quantization (PTQ).

- *Preparation*
- *Pre-trained model*
- *Quantize and export `.espd1`*
  - *Add test input/output*
  - *Quantized model inference & accuracy evaluation*

#### 3.1.1 Preparation

*Install ESP\_PPQ*

#### 3.1.2 Pre-trained model

```
python sin_model.py
```

Run `sin_model.py` . This script trains a simple Pytorch model to fit the sin function in the range  $[0, 2\pi]$ . After training, the corresponding `.pth` weights will be saved and the ONNX model will be exported.

---

**Note:** ESP-PPQ provides two interfaces, `espd1_quantize_onnx` and `espd1_quantize_torch`, to support ONNX models and PyTorch models. Other deep learning frameworks, such as TensorFlow, PaddlePaddle, etc., need to be converted to ONNX first.

- Convert TensorFlow to ONNX `tf2onnx`
  - Convert TFLite to ONNX `tf2onnx`
  - Convert TFLite to TensorFlow `tf2onnx`
  - Convert PaddlePaddle to ONNX `paddle2onnx`
-

### 3.1.3 Quantize and export .espdl

Reference [quantize\\_torch\\_model.py](#) and [quantize\\_onnx\\_model.py](#), learn how to use the `espdl_quantize_onnx` and `espdl_quantize_torch` interfaces to quantize and export the `.espdl` model.

After executing the script, three files will be exported:

- `** .espdl`: ESPDL model binary file, which can be directly used for chip reasoning.
- `** .info`: ESPDL model text file, used to debug and determine whether the `.espdl` model is exported correctly. Contains model structure, quantized model weights, test input/output and other information.
- `** .json`: Quantization information file, used to save and load quantization information.

---

**Note:**

1. The `.espdl` models of different platforms cannot be mixed, otherwise the inference results will be inaccurate. The ROUND strategy used by ESP32S3 is `ROUND_HALF_UP`, and the one used by ESP32P4 is `ROUND_HALF_EVEN`.
  2. The quantization strategy currently used by ESP-DL is symmetric quantization + POWER OF TWO.
- 

#### Add test input/output

To verify whether the inference results of the model on the board are correct, you first need to record a set of test input/output on the PC. By turning on the `export_test_values` option in the `api`, a set of test input/output can be saved in the `.espdl` model. One of the `input_shape` and `inputs` parameters must be specified. The `input_shape` parameter uses a random test input, while `inputs` can use a specific test input. The values of the test input/output can be viewed in the `.info` file. Search for `test inputs value` and `test outputs value` to view them.

#### Quantized model inference & accuracy evaluation

`espdl_quantize_onnx` and `espdl_quantize_torch` APIs will return `BaseGraph`. Use `BaseGraph` to build the corresponding `TorchExecutor` to use the quantized model for inference on the PC side.

```
executor = TorchExecutor(graph=quanted_graph, device=device)
output = executor(input)
```

The output obtained by quantized model inference can be used to calculate various accuracy metrics. Since the board-side `esp-dl` inference result can be aligned with `esp-ppq`, these metrics can be used directly to evaluate the accuracy of the quantized model.

---

**Note:**

1. Currently `esp-dl` only supports `batch_size` of 1, and does not support multi-batch or dynamic batch.
  2. The test input/output and the quantized model weights in the `.info` file are all 16-byte aligned. If the length is less than 16 bytes, it will be padded with 0.
- 

## 3.2 How to load & test & profile model

In this tutorial, we will show you how to load, test, profile an `espdl` model. [example](#)

- *Preparation*
- *Load model from rodata*
- *Load model from partition*
- *Load model from sdcard*
- *Test whether on-board model inference is correct*
- *Profile model memory usage*
- *Profile model inference latency*

### 3.2.1 Preparation

1. *Install ESP\_IDF*
2. *how\_to\_quantize\_model*

### 3.2.2 Load model from rodata

1. **Add model file in CMakeLists.txt**

If you want to put the .espdl model file into the .rodata section of the FLASH chip, you need to add the following code in CMakeLists.txt. The first few lines should be placed before `idf_component_register()` and the last line should be placed after `idf_component_register()`.

```
idf_build_get_property(component_targets __COMPONENT_TARGETS)
if ("__idf_espessif_esp-dl" IN_LIST component_targets)
idf_component_get_property(espdl_dir espressif_esp-dl COMPONENT_DIR)
elseif("__idf_esp-dl" IN_LIST component_targets)
idf_component_get_property(espdl_dir esp-dl COMPONENT_DIR)
endif()
set(cmake_dir ${espdl_dir}/fbs_loader/cmake)
include(${cmake_dir}/utilities.cmake)
set(embed_files your_model_path/model_name.espdl) idf_component_register(...)

target_add_aligned_binary_data(${COMPONENT_LIB} ${embed_files} BINARY)
```

2. **Load the model in the program**

```
// "_binary_model_espdl_start" is composed of three parts: the prefix "binary",
↪ the filename "model_espdl", and the suffix "_start".
extern const uint8_t model_espdl[] asm("_binary_model_espdl_start");

dl::Model *model = new dl::Model((const char *)model_espdl, fbs::MODEL_
↪LOCATION_IN_FLASH_RODATA);

// Keep parameter in FLASH, saves PSRAM/internal RAM, lower performance.
// dl::Model *model = new dl::Model((const char *)model_espdl, fbs::MODEL_
↪LOCATION_IN_FLASH_RODATA, 0, dl::MEMORY_MANAGER_GREEDY, nullptr, false);
```

#### Note:

1. When using *Load model from rodata*, since the .rodata section belongs to the app partition, the model file will be flashed every time the code is modified. If the model file is large, you may need to adjust the size of the app partition. Using *Load model from partition* or *Load model from sdcard* can avoid repeatedly flashing the model, which helps to reduce the flashing time.
2. When using *Load model from rodata* or *Load model from partition*, turning off the `param_copy` option in the Model constructor can avoid copying the model weights in PSRAM or internal RAM to FLASH. This can reduce the use of PSRAM or internal RAM. However, since the frequency of PSRAM or internal RAM is higher than FLASH, the inference performance of the model will decrease.

### 3.2.3 Load model from partition

#### 1. Add model information in partition.csv

About `partition.csv`, please refer to the [partition table documentation](#).

```
# Name, Type, SubType, Offset, Size, Flags
factory, app, factory, 0x010000, 4000K,
model, data, spiffs, , 4000K,
```

The Name field of the model can be any meaningful name, but cannot exceed 16 bytes, including a null byte (the content after that will be truncated). The SubType field must be `spiffs`. The Offset can be left blank after other partitions and will be automatically calculated. Size must be larger than the size of the model file.

#### 2. Add model flashing information in CMakeLists.txt

```
idf_component_register(...)
set(image_file your_model_path)
esptool_py_flash_to_partition(flash "model" "${image_file}")
```

The second parameter in `esptool_py_flash_to_partition` must be consistent with the Name field in `partition.csv`.

#### 3. Load the model in the program

```
dl::Model *model = new dl::Model("model", fbs::MODEL_LOCATION_IN_FLASH_PARTITION);

// Keep parameter in flash, saves PSRAM/internal RAM, lower performance.
// dl::Model *model = new dl::Model("model", fbs::MODEL_LOCATION_IN_FLASH_
↳PARTITION, 0, dl::MEMORY_MANAGER_GREEDY,
// nullptr, false);
```

The first parameter of the constructor must be consistent with the Name field in `partition.csv`.

---

**Note:** Use `idf.py app-flash` instead of `idf.py flash` to flash only the app partition without flashing the model partition, this can reduce the flashing time.

---

### 3.2.4 Load model from sdcard

#### 1. Check if sdcard is in the correct format

Back up the data first, then try to mount it on the board. If the sdcard is not in the correct format, it will be automatically formatted to the correct format.

- If using [BSP\(Board Support Package\)](#)

Enable `CONFIG_BSP_SD_FORMAT_ON_MOUNT_FAIL` option in `menuconfig`.

```
ESP_ERROR_CHECK(bsp_sdcard_mount());
```

- If not using [BSP\(Board Support Package\)](#)

Set `format_if_mount_failed` in `esp_vfs_fat_sdmmc_mount_config_t` structure to `true`.

```
esp_vfs_fat_sdmmc_mount_config_t mount_config = {
    .format_if_mount_failed = true,
    .max_files = 5,
    .allocation_unit_size = 16 * 1024
};
// Mount sdcard.
```

#### 2. Copy model to sdcard

Copy `.espd` model to sdcard.

#### 3. Load the model in the program

- If using BSP(Board Support Package)

```
ESP_ERROR_CHECK(bsp_sdcard_mount());
const char *model_path = "/your_sdcard_mount_point/your_model_path/model_name.
↳espd1";
Model *model = new Model(model_path, fbs::MODEL_LOCATION_IN_SDCARD);
```

- If not using BSP(Board Support Package)

```
// Mount sdcard.
const char *model_path = "/your_sdcard_mount_point/your_model_path/model_name.
↳espd1";
Model *model = new Model(model_path, fbs::MODEL_LOCATION_IN_SDCARD);
```

**Note:** When using *load model from sdcard*, the model loading process will take longer because the model data needs to be copied from sdcard to PSRAM or internal RAM. This method is useful if your FLASH is tight.

### 3.2.5 Test whether on-board model inference is correct

In order to test whether on-board model inference is correct, the `.espd1` model needs to *add test input/output* when exporting. When actually deploying, you can export a version without test input and output to reduce the size of the model file.

```
ESP_ERROR_CHECK(model->test());
```

### 3.2.6 Profile model memory usage

```
model->profile_memory();
```

| Name         | Explanation   |
|--------------|---|
| fbs_model    | Flatbuffers model size. Contains model parameters, test input/output, model parameter shapes, model structure etc.                          |
| param_in_fm  | Model parameter size (inside flatbuffers model).  |
| param_out_fm | Model parameter size (outside flatbuffers model).   |
| mem_manager  | Memory size allocated by the memory manager. Model input/output and intermediate calculation results will use this space.                   |
| others       | Space required for class member variables, extra part for alignment during heap_caps_aligned_alloc / heap_caps_aligned_calloc (very small). |

### 3.2.7 Profile model inference latency

```
model->profile_module();
```

By default, the model modules are printed in ONNX topological sort. If you want to sort by the latency of each module, you can set the input parameter of `profile_module` to `True`.

```
model->profile_module(true);
```

## 3.3 How to run model



In this tutorial, we will introduce the most basic model inference process. [example](#)

- *Preparation*
- *Load model*
- *Get model input/output.*
- *Quantize Input*
  - *Quantize a single value*
  - *Quantize `dl::TensorBase`*
- *Dequantize output*
  - *Dequantize a single value*
  - *Dequantize `dl::TensorBase`*
- *Model Inference*

### 3.3.1 Preparation

*Install ESP\_IDF*

### 3.3.2 Load model

*How to load model*

### 3.3.3 Get model input/output.

```
std::map<std::string, dl::TensorBase *> model_inputs = model->get_inputs();
dl::TensorBase *model_input = model_inputs.begin()->second;
std::map<std::string, dl::TensorBase *> model_outputs = model->get_outputs();
dl::TensorBase *model_output = model_outputs.begin()->second;
```

You can get the input/output names and the corresponding `dl::TensorBase` with `get_inputs()` and `get_outputs()` api. For more information, see [dl::TensorBase documentation](#).

**Note:** ESP-DL's memory manager allocates a whole block of memory for each model's input/intermediate result/output. Since they share this memory, when the model is inferencing, the later results will overwrite the previous results. In other words, the data in `model_input` may be overwritten by `model_output` or other intermediate results after the model inference is completed.

### 3.3.4 Quantize Input

8-bit and 16-bit quantized models accept inputs of type `int8_t` and `int16_t` respectively. `float` inputs must be quantized to the one of them according to `exponent` before being fed into the model. Calculation formula:

$$Q = \text{Clip} \left( \text{Round} \left( \frac{R}{\text{Scale}} \right), \text{MIN}, \text{MAX} \right)$$

$$\text{Scale} = 2^{\text{Exp}}$$

Where:

- R is the floating point number to be quantized.
- Q is the integer value after quantization, which needs to be clipped within the range [MIN, MAX].
- MIN is the minimum integer value, when 8bit, MIN = -128, when 16bit, MIN = -32768.
- MAX is the maximum integer value, when 8bit, MAX = 127, when 16bit, MAX = 32767.

### Quantize a single value

```
float input_v = VALUE;
// Note that dl::quantize accepts inverse of scale as the second input, so we use
↳DL_RESCALE here.
int8_t quant_input_v = dl::quantize<int8_t>(input_v, DL_RESCALE(model_input->
↳exponent));
```

### Quantize dl::TensorBase

```
// assume that input_tensor already contains the float input data.
dl::TensorBase *input_tensor;
model_input->assign(input_tensor);
```

### 3.3.5 Dequantize output

8bit and 16bit quantized model, get `int8_t` and `int16_t` type output respectively. Must be dequantized according to exponent to get floating point output. Calculation formula:

$$R' = Q \times \text{Scale}$$

$$\text{Scale} = 2^{\text{Exp}}$$

Where:

- $R'$  is the approximate floating point value recovered after dequantization.
- $Q$  is the integer value after quantization.

### Dequantize a single value

```
int8_t quant_output_v = VALUE;
float output_v = dl::dequantize(quant_output_v, DL_SCALE(model_output->exponent));
```

### Dequantize dl::TensorBase

```
// create a TensorBase filled with 0 of shape [1, 1]
dl::TensorBase *output_tensor = new dl::TensorBase({1, 1}, nullptr, 0, dl::DATA_
↳TYPE_FLOAT);
output_tensor->assign(model_output);
```

### 3.3.6 Model Inference

See:

- [example](#)
- `void dl::Model::run(runtime_mode_t mode)`
- `void dl::Model::run(TensorBase *input, runtime_mode_t mode)`
- `void dl::Model::run(std::map<std::string, TensorBase*> &user_inputs, runtime_mode_t mode, std::map<std::string, TensorBase*> user_outputs)`

## 3.4 How to deploy MobileNetV2

In this tutorial, we will introduce how to quantize a pre-trained MobileNetV2 model using ESP-PPQ and deploy the quantized MobileNetV2 model using ESP-DL.

- *Preparation*
- *Model quantization*
  - *Pre-trained model*
  - *Calibration dataset*
  - *8bit default configuration quantization*
  - *Mixed precision quantization*
  - *Layerwise equalization quantization*
- *Model deployment*
  - *Image classification base class*
  - *Pre-process*
  - *Post-process*

### 3.4.1 Preparation

1. *Install ESP\_IDF*
2. *Install ESP\_PPQ*

### 3.4.2 Model quantization

Quantization script

#### Pre-trained model

Load the pre-trained model of MobileNet\_v2 from torchvision. You can also download it from [ONNX models](#) or [TensorFlow models](#):

```
import torchvision
from torchvision.models.mobilenetv2 import MobileNet_V2_Weights

model = torchvision.models.mobilenet.mobilenet_v2(weights=MobileNet_V2_Weights.
↳ IMAGENET1K_V1)
```

#### Calibration dataset

The calibration dataset needs to be consistent with your model input format. The calibration dataset needs to cover all possible situations of your model input as much as possible to better quantize the model. Here we take the ImageNet dataset as an example to demonstrate how to prepare the calibration dataset.

Use torchvision to load the ImageNet dataset:

```
import torchvision.datasets as datasets
from torch.utils.data.dataset import Subset

dataset = datasets.ImageFolder(
    CALIB_DIR,
    transforms.Compose(
        [
            transforms.Resize(256),
```

(continues on next page)

(continued from previous page)

```

        transforms.CenterCrop(224),
        transforms.ToTensor(),
        transforms.Normalize(
            mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]
        ),
    ],
),
)
dataset = Subset(dataset, indices=[_ for _ in range(0, 1024)])
dataloader = DataLoader(
    dataset=dataset,
    batch_size=BATCH_SIZE,
    shuffle=False,
    num_workers=4,
    pin_memory=False,
    collate_fn=collate_fn1,
)
    
```

### 8bit default configuration quantization

#### Quantization settings

```

























target="esp32p4"
num_of_bits=8
batch_size=32
quant_setting = QuantizationSettingFactory.espdn_setting() # default setting
    
```

#### Quantization results

| Layer  | NOISE:SIGNAL POWER RATIO |
|--|--------------------------|
| /features/features.16/conv/conv.2/Conv:          | 48.831%                  |
| /features/features.15/conv/conv.2/Conv:          | 45.268%                  |
| /features/features.17/conv/conv.2/Conv:          | 43.112%                  |
| /features/features.18/features.18.0/Conv:        | 41.586%                  |
| /features/features.14/conv/conv.2/Conv:          | 41.135%                  |
| /features/features.13/conv/conv.2/Conv:          | 35.090%                  |
| /features/features.17/conv/conv.0/conv.0.0/Conv: | 32.895%                  |
| /features/features.16/conv/conv.1/conv.1.0/Conv: | 29.226%                  |
| /features/features.12/conv/conv.2/Conv:          | 28.895%                  |
| /features/features.16/conv/conv.0/conv.0.0/Conv: | 27.808%                  |
| /features/features.7/conv/conv.2/Conv:           | 27.675%                  |
| /features/features.10/conv/conv.2/Conv:          | 26.292%                  |
| /features/features.11/conv/conv.2/Conv:          | 26.085%                  |
| /features/features.6/conv/conv.2/Conv:           | 25.892%                  |
| /classifier/classifier.1/Gemm:                   | 25.591%                  |
| /features/features.15/conv/conv.0/conv.0.0/Conv: | 25.323%                  |
| /features/features.4/conv/conv.2/Conv:           | 24.787%                  |
| /features/features.15/conv/conv.1/conv.1.0/Conv: | 24.354%                  |
| /features/features.14/conv/conv.1/conv.1.0/Conv: | 20.207%                  |
| /features/features.9/conv/conv.2/Conv:           | 19.808%                  |
| /features/features.14/conv/conv.0/conv.0.0/Conv: | 18.465%                  |
| /features/features.5/conv/conv.2/Conv:           | 17.868%                  |
| /features/features.12/conv/conv.1/conv.1.0/Conv: | 16.589%                  |
| /features/features.13/conv/conv.1/conv.1.0/Conv: | 16.143%                  |
| /features/features.11/conv/conv.1/conv.1.0/Conv: | 15.382%                  |
| /features/features.3/conv/conv.2/Conv:           | 15.105%                  |
| /features/features.13/conv/conv.0/conv.0.0/Conv: | 15.029%                  |
| /features/features.10/conv/conv.1/conv.1.0/Conv: | 14.875%                  |
| /features/features.2/conv/conv.2/Conv:           | 14.869%                  |

(continues on next page)


































(continued from previous page)

|  |  |         |
|--|--|---------|
| /features/features.11/conv/conv.0/conv.0.0/Conv: |  | 14.552% |
| /features/features.9/conv/conv.1/conv.1.0/Conv:  |  | 14.050% |
| /features/features.8/conv/conv.1/conv.1.0/Conv:  |  | 13.929% |
| /features/features.8/conv/conv.2/Conv:           |  | 13.833% |
| /features/features.12/conv/conv.0/conv.0.0/Conv: |  | 13.684% |
| /features/features.7/conv/conv.0/conv.0.0/Conv:  |  | 12.942% |
| /features/features.6/conv/conv.1/conv.1.0/Conv:  |  | 12.765% |
| /features/features.10/conv/conv.0/conv.0.0/Conv: |  | 12.251% |
| /features/features.5/conv/conv.1/conv.1.0/Conv:  |  | 11.186% |
| /features/features.17/conv/conv.1/conv.1.0/Conv: |  | 11.070% |
| /features/features.9/conv/conv.0/conv.0.0/Conv:  |  | 10.371% |
| /features/features.4/conv/conv.1/conv.1.0/Conv:  |  | 10.356% |
| /features/features.6/conv/conv.0/conv.0.0/Conv:  |  | 10.149% |
| /features/features.4/conv/conv.0/conv.0.0/Conv:  |  | 9.472%  |
| /features/features.8/conv/conv.0/conv.0.0/Conv:  |  | 9.232%  |
| /features/features.3/conv/conv.1/conv.1.0/Conv:  |  | 9.187%  |
| /features/features.1/conv/conv.1/Conv:           |  | 8.770%  |
| /features/features.5/conv/conv.0/conv.0.0/Conv:  |  | 8.408%  |
| /features/features.7/conv/conv.1/conv.1.0/Conv:  |  | 8.151%  |
| /features/features.2/conv/conv.1/conv.1.0/Conv:  |  | 7.156%  |
| /features/features.3/conv/conv.0/conv.0.0/Conv:  |  | 6.328%  |
| /features/features.2/conv/conv.0/conv.0.0/Conv:  |  | 5.392%  |
| /features/features.1/conv/conv.0/conv.0.0/Conv:  |  | 0.875%  |
| /features/features.0/features.0.0/Conv:          |  | 0.119%  |

Analysing Layerwise quantization error:: 100

↪% | 

↪53/53 [08:44<00:00, 9.91s/it]

| Layer  |  | NOISE: SIGNAL POWER RATIO |
|--|--|---------------------------|
| /features/features.1/conv/conv.0/conv.0.0/Conv:  |  | 14.303%                   |
| /features/features.0/features.0.0/Conv:          |   | 0.844%                    |
| /features/features.1/conv/conv.1/Conv:           |   | 0.667%                    |
| /features/features.2/conv/conv.1/conv.1.0/Conv:  |   | 0.574%                    |
| /features/features.3/conv/conv.1/conv.1.0/Conv:  |   | 0.419%                    |
| /features/features.15/conv/conv.1/conv.1.0/Conv: |   | 0.272%                    |
| /features/features.9/conv/conv.1/conv.1.0/Conv:  |   | 0.238%                    |
| /features/features.17/conv/conv.1/conv.1.0/Conv: |   | 0.214%                    |
| /features/features.4/conv/conv.1/conv.1.0/Conv:  |   | 0.180%                    |
| /features/features.11/conv/conv.1/conv.1.0/Conv: |   | 0.151%                    |
| /features/features.12/conv/conv.1/conv.1.0/Conv: |   | 0.148%                    |
| /features/features.16/conv/conv.1/conv.1.0/Conv: |   | 0.146%                    |
| /features/features.14/conv/conv.2/Conv:          |   | 0.136%                    |
| /features/features.13/conv/conv.1/conv.1.0/Conv: |   | 0.105%                    |
| /features/features.6/conv/conv.1/conv.1.0/Conv:  |   | 0.105%                    |
| /features/features.8/conv/conv.1/conv.1.0/Conv:  |   | 0.083%                    |
| /features/features.7/conv/conv.2/Conv:           |   | 0.076%                    |
| /features/features.5/conv/conv.1/conv.1.0/Conv:  |   | 0.076%                    |
| /features/features.3/conv/conv.2/Conv:           |   | 0.075%                    |
| /features/features.16/conv/conv.2/Conv:          |   | 0.074%                    |
| /features/features.13/conv/conv.0/conv.0.0/Conv: |   | 0.072%                    |
| /features/features.15/conv/conv.2/Conv:          |   | 0.066%                    |
| /features/features.4/conv/conv.2/Conv:           |   | 0.065%                    |
| /features/features.11/conv/conv.2/Conv:          |   | 0.063%                    |
| /classifier/classifier.1/Gemm:                   |   | 0.063%                    |
| /features/features.2/conv/conv.0/conv.0.0/Conv:  |   | 0.054%                    |
| /features/features.13/conv/conv.2/Conv:          |   | 0.050%                    |
| /features/features.10/conv/conv.1/conv.1.0/Conv: |   | 0.042%                    |
| /features/features.17/conv/conv.0/conv.0.0/Conv: |   | 0.040%                    |
| /features/features.2/conv/conv.2/Conv:           |   | 0.038%                    |
| /features/features.4/conv/conv.0/conv.0.0/Conv:  |   | 0.034%                    |
| /features/features.17/conv/conv.2/Conv:          |   | 0.030%                    |
| /features/features.14/conv/conv.0/conv.0.0/Conv: |   | 0.025%                    |

(continues on next page)

(continued from previous page)

```

/features/features.16/conv/conv.0/conv.0.0/Conv: | 0.024%
/features/features.10/conv/conv.2/Conv: | 0.022%
/features/features.11/conv/conv.0/conv.0.0/Conv: | 0.021%
/features/features.9/conv/conv.2/Conv: | 0.021%
/features/features.14/conv/conv.1/conv.1.0/Conv: | 0.020%
/features/features.7/conv/conv.1/conv.1.0/Conv: | 0.020%
/features/features.5/conv/conv.2/Conv: | 0.019%
/features/features.8/conv/conv.2/Conv: | 0.018%
/features/features.12/conv/conv.2/Conv: | 0.017%
/features/features.6/conv/conv.2/Conv: | 0.014%
/features/features.7/conv/conv.0/conv.0.0/Conv: | 0.014%
/features/features.3/conv/conv.0/conv.0.0/Conv: | 0.013%
/features/features.12/conv/conv.0/conv.0.0/Conv: | 0.009%
/features/features.15/conv/conv.0/conv.0.0/Conv: | 0.008%
/features/features.5/conv/conv.0/conv.0.0/Conv: | 0.006%
/features/features.6/conv/conv.0/conv.0.0/Conv: | 0.005%
/features/features.9/conv/conv.0/conv.0.0/Conv: | 0.003%
/features/features.18/features.18.0/Conv: | 0.002%
/features/features.10/conv/conv.0/conv.0.0/Conv: | 0.002%
/features/features.8/conv/conv.0/conv.0.0/Conv: | 0.002%
* Prec@1 60.500 Prec@5 83.275*

```

### Quantization error analysis

The top1 accuracy after quantization is only 60.5%, which is far from the accuracy of the float model (71.878%). The quantization model has a large loss in accuracy, including:

- **Graphwise Error**

The last layer of the model is `/classifier/classifier.1/Gemm`, and the cumulative error of this layer is 25.591%. In experience, the cumulative error of the last layer is less than 10%, and the accuracy loss of the quantization model is small.

- **Layerwise error**

Observing the Layerwise error, it is found that the errors of most layers are less than 1%, indicating that the quantization errors of most layers are small, and only a few layers have large errors. We can choose to quantize the layers with large errors using `int16`. For details, please see mixed precision quantization.

### Mixed precision quantization

#### Quantization settings

```

from ppq.api import get_target_platform
target="esp32p4"
num_of_bits=8
batch_size=32

# The following layers are quantized using int16
quant_setting = QuantizationSettingFactory.espd1_setting()
quant_setting.dispatching_table.append("/features/features.1/conv/conv.0/conv.0.0/
↳Conv", get_target_platform(TARGET, 16))
quant_setting.dispatching_table.append("/features/features.1/conv/conv.0/conv.0.2/
↳Clip", get_target_platform(TARGET, 16))

```

#### Quantization results

| Layer  | NOISE:SIGNAL POWER RATIO |
|--|--------------------------|
| /features/features.16/conv/conv.2/Conv:          | 31.585%                  |
| /features/features.15/conv/conv.2/Conv:          | 29.253%                  |
| /features/features.17/conv/conv.0/conv.0.0/Conv: | 25.077%                  |
| /features/features.14/conv/conv.2/Conv:          | 24.819%                  |

(continues on next page)

(continued from previous page)

|  |            |         |
|--|------------|---------|
| /features/features.17/conv/conv.2/Conv:          | ██████████ | 19.546% |
| /features/features.13/conv/conv.2/Conv:          | ██████████ | 19.283% |
| /features/features.16/conv/conv.0/conv.0.0/Conv: | ██████████ | 18.764% |
| /features/features.16/conv/conv.1/conv.1.0/Conv: | ██████████ | 18.596% |
| /features/features.18/features.18.0/Conv:        | ██████████ | 18.541% |
| /features/features.15/conv/conv.0/conv.0.0/Conv: | ██████████ | 15.633% |
| /features/features.12/conv/conv.2/Conv:          | ██████████ | 14.784% |
| /features/features.15/conv/conv.1/conv.1.0/Conv: | ██████████ | 14.773% |
| /features/features.14/conv/conv.1/conv.1.0/Conv: | ██████████ | 13.700% |
| /features/features.6/conv/conv.2/Conv:           | ██████████ | 12.824% |
| /features/features.10/conv/conv.2/Conv:          | ██████████ | 11.727% |
| /features/features.14/conv/conv.0/conv.0.0/Conv: | ██████████ | 10.612% |
| /features/features.11/conv/conv.2/Conv:          | ██████████ | 10.262% |
| /features/features.9/conv/conv.2/Conv:           | ██████████ | 9.967%  |
| /classifier/classifier.1/Gemm:                   | ██████████ | 9.117%  |
| /features/features.5/conv/conv.2/Conv:           | ██████████ | 8.915%  |
| /features/features.7/conv/conv.2/Conv:           | ██████████ | 8.690%  |
| /features/features.3/conv/conv.2/Conv:           | ██████████ | 8.586%  |
| /features/features.4/conv/conv.2/Conv:           | ██████████ | 7.525%  |
| /features/features.13/conv/conv.1/conv.1.0/Conv: | ██████████ | 7.432%  |
| /features/features.12/conv/conv.1/conv.1.0/Conv: | ██████████ | 7.317%  |
| /features/features.13/conv/conv.0/conv.0.0/Conv: | ██████████ | 6.848%  |
| /features/features.8/conv/conv.2/Conv:           | ██████████ | 6.711%  |
| /features/features.10/conv/conv.1/conv.1.0/Conv: | ██████████ | 6.100%  |
| /features/features.8/conv/conv.1/conv.1.0/Conv:  | ██████████ | 6.043%  |
| /features/features.11/conv/conv.1/conv.1.0/Conv: | ██████████ | 5.962%  |
| /features/features.9/conv/conv.1/conv.1.0/Conv:  | ██████████ | 5.873%  |
| /features/features.12/conv/conv.0/conv.0.0/Conv: | ██████████ | 5.833%  |
| /features/features.7/conv/conv.0/conv.0.0/Conv:  | ██████████ | 5.832%  |
| /features/features.11/conv/conv.0/conv.0.0/Conv: | ██████████ | 5.736%  |
| /features/features.6/conv/conv.1/conv.1.0/Conv:  | ██████████ | 5.639%  |
| /features/features.5/conv/conv.1/conv.1.0/Conv:  | ██████████ | 5.017%  |
| /features/features.10/conv/conv.0/conv.0.0/Conv: | ██████████ | 4.963%  |
| /features/features.17/conv/conv.1/conv.1.0/Conv: | ██████████ | 4.870%  |
| /features/features.3/conv/conv.1/conv.1.0/Conv:  | ██████████ | 4.655%  |
| /features/features.2/conv/conv.2/Conv:           | ██████████ | 4.650%  |
| /features/features.4/conv/conv.0/conv.0.0/Conv:  | ██████████ | 4.648%  |
| /features/features.1/conv/conv.1/Conv:           | ██████████ | 4.318%  |
| /features/features.9/conv/conv.0/conv.0.0/Conv:  | ██████████ | 3.849%  |
| /features/features.6/conv/conv.0/conv.0.0/Conv:  | ██████████ | 3.712%  |
| /features/features.4/conv/conv.1/conv.1.0/Conv:  | ██████████ | 3.394%  |
| /features/features.8/conv/conv.0/conv.0.0/Conv:  | ██████████ | 3.391%  |
| /features/features.7/conv/conv.1/conv.1.0/Conv:  | ██████████ | 2.713%  |
| /features/features.2/conv/conv.1/conv.1.0/Conv:  | ██████████ | 2.637%  |
| /features/features.2/conv/conv.0/conv.0.0/Conv:  | ██████████ | 2.602%  |
| /features/features.5/conv/conv.0/conv.0.0/Conv:  | ██████████ | 2.397%  |
| /features/features.3/conv/conv.0/conv.0.0/Conv:  | ██████████ | 1.759%  |
| /features/features.1/conv/conv.0/conv.0.0/Conv:  | ██████████ | 0.433%  |
| /features/features.0/features.0.0/Conv:          | ██████████ | 0.119%  |

Analysing Layerwise quantization error:: 100

↪ % | ██████████  
 ↪ 53/53 [08:27<00:00, 9.58s/it]

\*

| Layer  | NOISE:SIGNAL POWER RATIO |
|--|--------------------------|
| /features/features.1/conv/conv.1/Conv:           | ██████████   1.096%      |
| /features/features.0/features.0.0/Conv:          | ██████████   0.844%      |
| /features/features.2/conv/conv.1/conv.1.0/Conv:  | ██████████   0.574%      |
| /features/features.3/conv/conv.1/conv.1.0/Conv:  | ██████████   0.425%      |
| /features/features.15/conv/conv.1/conv.1.0/Conv: | ██████████   0.272%      |
| /features/features.9/conv/conv.1/conv.1.0/Conv:  | ██████████   0.238%      |
| /features/features.17/conv/conv.1/conv.1.0/Conv: | ██████████   0.214%      |

(continues on next page)

(continued from previous page)

|  |   |        |
|--|---|--------|
| /features/features.4/conv/conv.1/conv.1.0/Conv:  | ■ | 0.180% |
| /features/features.11/conv/conv.1/conv.1.0/Conv: | ■ | 0.151% |
| /features/features.12/conv/conv.1/conv.1.0/Conv: | ■ | 0.148% |
| /features/features.16/conv/conv.1/conv.1.0/Conv: | ■ | 0.146% |
| /features/features.14/conv/conv.2/Conv:          | ■ | 0.136% |
| /features/features.13/conv/conv.1/conv.1.0/Conv: | ■ | 0.105% |
| /features/features.6/conv/conv.1/conv.1.0/Conv:  | ■ | 0.105% |
| /features/features.8/conv/conv.1/conv.1.0/Conv:  | ■ | 0.083% |
| /features/features.5/conv/conv.1/conv.1.0/Conv:  | ■ | 0.076% |
| /features/features.3/conv/conv.2/Conv:           | ■ | 0.075% |
| /features/features.16/conv/conv.2/Conv:          | ■ | 0.074% |
| /features/features.13/conv/conv.0/conv.0.0/Conv: | ■ | 0.072% |
| /features/features.7/conv/conv.2/Conv:           | ■ | 0.071% |
| /features/features.15/conv/conv.2/Conv:          | ■ | 0.066% |
| /features/features.4/conv/conv.2/Conv:           | ■ | 0.065% |
| /features/features.11/conv/conv.2/Conv:          | ■ | 0.063% |
| /classifier/classifier.1/Gemm:                   | ■ | 0.063% |
| /features/features.13/conv/conv.2/Conv:          | ■ | 0.059% |
| /features/features.2/conv/conv.0/conv.0.0/Conv:  | ■ | 0.054% |
| /features/features.10/conv/conv.1/conv.1.0/Conv: | ■ | 0.042% |
| /features/features.17/conv/conv.0/conv.0.0/Conv: | ■ | 0.040% |
| /features/features.2/conv/conv.2/Conv:           | ■ | 0.038% |
| /features/features.4/conv/conv.0/conv.0.0/Conv:  | ■ | 0.034% |
| /features/features.17/conv/conv.2/Conv:          | ■ | 0.030% |
| /features/features.14/conv/conv.0/conv.0.0/Conv: | ■ | 0.025% |
| /features/features.16/conv/conv.0/conv.0.0/Conv: | ■ | 0.024% |
| /features/features.10/conv/conv.2/Conv:          | ■ | 0.022% |
| /features/features.11/conv/conv.0/conv.0.0/Conv: | ■ | 0.021% |
| /features/features.9/conv/conv.2/Conv:           | ■ | 0.021% |
| /features/features.14/conv/conv.1/conv.1.0/Conv: | ■ | 0.020% |
| /features/features.7/conv/conv.1/conv.1.0/Conv:  | ■ | 0.020% |
| /features/features.5/conv/conv.2/Conv:           | ■ | 0.019% |
| /features/features.8/conv/conv.2/Conv:           | ■ | 0.018% |
| /features/features.12/conv/conv.2/Conv:          | ■ | 0.017% |
| /features/features.1/conv/conv.0/conv.0.0/Conv:  | ■ | 0.017% |
| /features/features.6/conv/conv.2/Conv:           | ■ | 0.014% |
| /features/features.7/conv/conv.0/conv.0.0/Conv:  | ■ | 0.014% |
| /features/features.3/conv/conv.0/conv.0.0/Conv:  | ■ | 0.013% |
| /features/features.12/conv/conv.0/conv.0.0/Conv: | ■ | 0.009% |
| /features/features.15/conv/conv.0/conv.0.0/Conv: | ■ | 0.008% |
| /features/features.5/conv/conv.0/conv.0.0/Conv:  | ■ | 0.006% |
| /features/features.6/conv/conv.0/conv.0.0/Conv:  | ■ | 0.005% |
| /features/features.9/conv/conv.0/conv.0.0/Conv:  | ■ | 0.003% |
| /features/features.18/features.18.0/Conv:        | ■ | 0.002% |
| /features/features.10/conv/conv.0/conv.0.0/Conv: | ■ | 0.002% |
| /features/features.8/conv/conv.0/conv.0.0/Conv:  | ■ | 0.002% |

\* Prec@1 69.550 Prec@5 88.450\*

### Quantization Error Analysis

After replacing the layer with the largest error with 16-bit quantization, it can be observed that the model accuracy is significantly improved. The top1 accuracy after quantization is 69.550%, which is close to the accuracy of the float model (71.878%). The cumulative error of the last layer of the model `/classifier/classifier.1/Gemm` is 9.117%.

### Layerwise equalization quantization

This method is proposed in the paper [Data-Free Quantization Through Weight Equalization and Bias Correction](#). When using this method, the original ReLU6 in the MobilenetV2 model needs to be replaced with ReLU.



## Quantization Settings

```

import torch.nn as nn
def convert_relu6_to_relu(model):
    for child_name, child in model.named_children():
        if isinstance(child, nn.ReLU6):
            setattr(model, child_name, nn.ReLU())
        else:
            convert_relu6_to_relu(child)
    return model

# Replace ReLU6 with ReLU
model = convert_relu6_to_relu(model)
# Use layerwise equalization
quant_setting = QuantizationSettingFactory.espd1_setting()
quant_setting.equalization = True
quant_setting.equalization_setting.iterations = 4
quant_setting.equalization_setting.value_threshold = .4
quant_setting.equalization_setting.opt_level = 2
quant_setting.equalization_setting.interested_layers = None

```

| Layer  | NOISE:SIGNAL POWER RATIO |
|--|--------------------------|
| /features/features.16/conv/conv.2/Conv:          | 34.497%                  |
| /features/features.15/conv/conv.2/Conv:          | 30.813%                  |
| /features/features.14/conv/conv.2/Conv:          | 25.876%                  |
| /features/features.17/conv/conv.0/conv.0.0/Conv: | 24.498%                  |
| /features/features.17/conv/conv.2/Conv:          | 20.290%                  |
| /features/features.13/conv/conv.2/Conv:          | 20.177%                  |
| /features/features.16/conv/conv.0/conv.0.0/Conv: | 19.993%                  |
| /features/features.18/features.18.0/Conv:        | 19.536%                  |
| /features/features.16/conv/conv.1/conv.1.0/Conv: | 17.879%                  |
| /features/features.12/conv/conv.2/Conv:          | 17.150%                  |
| /features/features.15/conv/conv.0/conv.0.0/Conv: | 15.970%                  |
| /features/features.15/conv/conv.1/conv.1.0/Conv: | 15.254%                  |
| /features/features.1/conv/conv.1/Conv:           | 15.122%                  |
| /features/features.10/conv/conv.2/Conv:          | 14.917%                  |
| /features/features.6/conv/conv.2/Conv:           | 13.446%                  |
| /features/features.11/conv/conv.2/Conv:          | 12.533%                  |
| /features/features.9/conv/conv.2/Conv:           | 11.479%                  |
| /features/features.14/conv/conv.1/conv.1.0/Conv: | 11.470%                  |
| /features/features.5/conv/conv.2/Conv:           | 10.669%                  |
| /features/features.3/conv/conv.2/Conv:           | 10.526%                  |
| /features/features.14/conv/conv.0/conv.0.0/Conv: | 9.529%                   |
| /features/features.7/conv/conv.2/Conv:           | 9.500%                   |
| /classifier/classifier.1/Gemm:                   | 8.965%                   |
| /features/features.4/conv/conv.2/Conv:           | 8.674%                   |
| /features/features.12/conv/conv.1/conv.1.0/Conv: | 8.349%                   |
| /features/features.13/conv/conv.1/conv.1.0/Conv: | 8.068%                   |
| /features/features.8/conv/conv.2/Conv:           | 7.961%                   |
| /features/features.13/conv/conv.0/conv.0.0/Conv: | 7.451%                   |
| /features/features.10/conv/conv.1/conv.1.0/Conv: | 6.714%                   |
| /features/features.9/conv/conv.1/conv.1.0/Conv:  | 6.399%                   |
| /features/features.8/conv/conv.1/conv.1.0/Conv:  | 6.369%                   |
| /features/features.11/conv/conv.1/conv.1.0/Conv: | 6.222%                   |
| /features/features.2/conv/conv.2/Conv:           | 5.867%                   |
| /features/features.5/conv/conv.1/conv.1.0/Conv:  | 5.719%                   |
| /features/features.12/conv/conv.0/conv.0.0/Conv: | 5.546%                   |
| /features/features.6/conv/conv.1/conv.1.0/Conv:  | 5.414%                   |
| /features/features.10/conv/conv.0/conv.0.0/Conv: | 5.093%                   |
| /features/features.17/conv/conv.1/conv.1.0/Conv: | 4.951%                   |
| /features/features.11/conv/conv.0/conv.0.0/Conv: | 4.941%                   |
| /features/features.2/conv/conv.1/conv.1.0/Conv:  | 4.825%                   |

(continues on next page)

(continued from previous page)

|  |  |                          |
|--|--|--------------------------|
| /features/features.7/conv/conv.0/conv.0.0/Conv:  | ███                                      | 4.330%                   |
| /features/features.2/conv/conv.0/conv.0.0/Conv:  | ███                                      | 4.299%                   |
| /features/features.3/conv/conv.1/conv.1.0/Conv:  | ███                                      | 4.283%                   |
| /features/features.4/conv/conv.0/conv.0.0/Conv:  | ███                                      | 3.477%                   |
| /features/features.4/conv/conv.1/conv.1.0/Conv:  | ███                                      | 3.287%                   |
| /features/features.8/conv/conv.0/conv.0.0/Conv:  | ███                                      | 2.787%                   |
| /features/features.9/conv/conv.0/conv.0.0/Conv:  | ███                                      | 2.774%                   |
| /features/features.6/conv/conv.0/conv.0.0/Conv:  | ███                                      | 2.705%                   |
| /features/features.7/conv/conv.1/conv.1.0/Conv:  | ███                                      | 2.636%                   |
| /features/features.5/conv/conv.0/conv.0.0/Conv:  | ███                                      | 1.846%                   |
| /features/features.3/conv/conv.0/conv.0.0/Conv:  | ███                                      | 1.170%                   |
| /features/features.1/conv/conv.0/conv.0.0/Conv:  | ███                                      | 0.389%                   |
| /features/features.0/features.0.0/Conv:  | ███                                      | 0.025%                   |
| Analysing Layerwise quantization error:: 100%  ██████████   53/53 [07:46<00:00, 8.<br>↪80s/it] |  |                          |
| Layer  |  | NOISE:SIGNAL POWER RATIO |
| /features/features.1/conv/conv.0/conv.0.0/Conv:  | ██ | 0.989%                   |
| /features/features.0/features.0.0/Conv:  | ██████████████████████████████████████   | 0.845%                   |
| /features/features.16/conv/conv.2/Conv:  | ██████████████████████████████████       | 0.238%                   |
| /features/features.17/conv/conv.2/Conv:  | ██████████████████████████████           | 0.202%                   |
| /features/features.14/conv/conv.2/Conv:  | ██████████████████████████████           | 0.198%                   |
| /features/features.1/conv/conv.1/Conv:   | ██████████████████████████               | 0.192%                   |
| /features/features.15/conv/conv.2/Conv:  | ██████████████████████████               | 0.145%                   |
| /features/features.4/conv/conv.2/Conv:   | ██████████████████████                   | 0.120%                   |
| /features/features.2/conv/conv.2/Conv:   | ██████████████████████                   | 0.111%                   |
| /features/features.2/conv/conv.1/conv.1.0/Conv:  | ██████████████████████                   | 0.079%                   |
| /classifier/classifier.1/Gemm:   | ██████████████████                       | 0.062%                   |
| /features/features.13/conv/conv.2/Conv:  | ██████████████████                       | 0.050%                   |
| /features/features.3/conv/conv.2/Conv:   | ██████████████████                       | 0.050%                   |
| /features/features.12/conv/conv.2/Conv:  | ██████████████████                       | 0.050%                   |
| /features/features.5/conv/conv.1/conv.1.0/Conv:  | ██████████████████                       | 0.047%                   |
| /features/features.3/conv/conv.1/conv.1.0/Conv:  | ██████████████████                       | 0.046%                   |
| /features/features.7/conv/conv.2/Conv:   | ██████████████████                       | 0.045%                   |
| /features/features.5/conv/conv.2/Conv:   | ██████████████████                       | 0.030%                   |
| /features/features.11/conv/conv.2/Conv:  | ██████████████████                       | 0.028%                   |
| /features/features.6/conv/conv.2/Conv:   | ██████████████████                       | 0.027%                   |
| /features/features.6/conv/conv.1/conv.1.0/Conv:  | ██████████████████                       | 0.026%                   |
| /features/features.4/conv/conv.0/conv.0.0/Conv:  | ███                                      | 0.025%                   |
| /features/features.15/conv/conv.1/conv.1.0/Conv:   | ███                                      | 0.023%                   |
| /features/features.8/conv/conv.1/conv.1.0/Conv:  | ███                                      | 0.021%                   |
| /features/features.10/conv/conv.2/Conv:  | ███                                      | 0.020%                   |
| /features/features.11/conv/conv.1/conv.1.0/Conv:   | ███                                      | 0.020%                   |
| /features/features.16/conv/conv.1/conv.1.0/Conv:   | ███                                      | 0.017%                   |
| /features/features.14/conv/conv.0/conv.0.0/Conv:   | ███                                      | 0.016%                   |
| /features/features.4/conv/conv.1/conv.1.0/Conv:  | ███                                      | 0.012%                   |
| /features/features.13/conv/conv.1/conv.1.0/Conv:   | ███                                      | 0.012%                   |
| /features/features.13/conv/conv.0/conv.0.0/Conv:   | ███                                      | 0.012%                   |
| /features/features.12/conv/conv.1/conv.1.0/Conv:   | ███                                      | 0.012%                   |
| /features/features.17/conv/conv.0/conv.0.0/Conv:   | ███                                      | 0.011%                   |
| /features/features.12/conv/conv.0/conv.0.0/Conv:   | ███                                      | 0.011%                   |
| /features/features.2/conv/conv.0/conv.0.0/Conv:  | ███                                      | 0.010%                   |
| /features/features.9/conv/conv.2/Conv:   | ███                                      | 0.008%                   |
| /features/features.8/conv/conv.2/Conv:   | ███                                      | 0.008%                   |
| /features/features.10/conv/conv.1/conv.1.0/Conv:   | ███                                      | 0.008%                   |
| /features/features.16/conv/conv.0/conv.0.0/Conv:   | ███                                      | 0.008%                   |
| /features/features.7/conv/conv.0/conv.0.0/Conv:  | ███                                      | 0.008%                   |
| /features/features.10/conv/conv.0/conv.0.0/Conv:   | ███                                      | 0.006%                   |
| /features/features.15/conv/conv.0/conv.0.0/Conv:   | ███                                      | 0.005%                   |
| /features/features.3/conv/conv.0/conv.0.0/Conv:  | ███                                      | 0.004%                   |
| /features/features.11/conv/conv.0/conv.0.0/Conv:   | ███                                      | 0.004%                   |
| /features/features.18/features.18.0/Conv:  | ███                                      | 0.003%                   |

(continues on next page)

(continued from previous page)

```

/features/features.5/conv/conv.0/conv.0.0/Conv: | 0.003%
/features/features.9/conv/conv.1/conv.1.0/Conv: | 0.003%
/features/features.6/conv/conv.0/conv.0.0/Conv: | 0.003%
/features/features.7/conv/conv.1/conv.1.0/Conv: | 0.003%
/features/features.17/conv/conv.1/conv.1.0/Conv: | 0.002%
/features/features.14/conv/conv.1/conv.1.0/Conv: | 0.002%
/features/features.8/conv/conv.0/conv.0.0/Conv: | 0.001%
/features/features.9/conv/conv.0/conv.0.0/Conv: | 0.001%

* Prec@1 69.800 Prec@5 88.550

```

### Quantization Error Analysis

Note that applying layerwise equalization to 8-bit quantization helps reduce quantization loss. The cumulative error of the last layer of the model `/classifier/classifier.1/Gemm` is 8.965%. The top1 accuracy after quantization is 69.800%, which is closer to the accuracy of the float model (71.878%) and higher than the quantization accuracy of mixed precision quantization.

**Note:** To further reduce the quantization error, you can try using QAT (Quantization Aware Training). For specific methods, please refer to [PPQ QAT example](#).

## 3.4.3 Model deployment

examples

### Image classification base class

- [dl\\_cls\\_base.hpp](#)
- [dl\\_cls\\_base.cpp](#)

### Pre-process

`ImagePreprocessor` class contains the common pre-process pipeline, color conversion, crop, resize, normalization, quantize.

- [dl\\_image\\_preprocessor.hpp](#)
- [dl\\_image\\_preprocessor.cpp](#)

### Post-process

- [dl\\_cls\\_postprocessor.hpp](#)
- [dl\\_cls\\_postprocessor.cpp](#)
- [imagenet\\_cls\\_postprocessor.hpp](#)
- [imagenet\\_cls\\_postprocessor.cpp](#)

## 3.5 How to deploy YOLO11n

In this tutorial, we will introduce how to quantize a pre-trained YOLO11n model using ESP-PPQ and deploy the quantized YOLO11n model using ESP-DL.

- *Preparation*
- *Model quantization*
  - *Pre-trained Model*
  - *Calibration Dataset*
  - *8bit default configuration quantization*
  - *Mixed-Precision + Horizontal Layer Split Pass Quantization*
  - *Quantization-Aware Training*
- *Model deployment*
  - *Object detection base class*
  - *Pre-process*
  - *Post-process*

### 3.5.1 Preparation

1. 安装 *ESP\_IDF*
2. 安装 *ESP\_PPQ*

### 3.5.2 Model quantization

#### Pre-trained Model

You can download pre-trained yolo11n model from [Ultralytics release](#).

Currently, ESP-PPQ supports ONNX, PyTorch, and TensorFlow models. During the quantization process, PyTorch and TensorFlow models are first converted to ONNX models, so the pre-trained yolo11n model needs to be converted to an ONNX model.

Specifically, refer to the script `export_onnx.py` to convert the pre-trained yolo11n model to an ONNX model.

In the script, we have overridden the forward method of the Detect class, which offers following advantages:

- **Faster inference.** Compared to the original yolo11n model, operations related to decoding bounding boxes in Detect head are moved from the inference pass to the post-processing phase, resulting in a significant reduction in inference latency. On one hand, operations like `Conv`, `Transpose`, `Slice`, `Split` and `Concat` are time-consuming when applied during inference pass. On the other hand, the inference outputs are first filtered using a score threshold before decoding the boxes in the post-processing pass, which significantly reduces the number of calculations, thereby accelerating the overall inference speed.
- **Lower quantization Error.** The `Concat` and `Add` operators adopt joint quantization in ESP-PPQ. To reduce quantization errors, the box and score are output by separate branches, rather than being concatenated, due to the significant difference in their ranges. Similarly, since the ranges of the two inputs of `Add` and `Sub` differ significantly, the calculations are performed in the post-processing phase to avoid quantization errors.

#### Calibration Dataset

The calibration dataset needs to match the input format of the model. The calibration dataset should cover all possible input scenarios to better quantize the model. Here, the calibration dataset used in this example is `calib_yolo11n`.

#### 8bit default configuration quantization

##### Quantization settings

```
target="esp32p4"
num_of_bits=8
batch_size=32
quant_setting = QuantizationSettingFactory.espd1_setting() # default setting
```

## Quantization results

| Layer  | NOISE:SIGNAL POWER RATIO |
|--|--------------------------|
| /model.10/m/m.0/ffn/ffn.1/conv/Conv:         | 36.042%                  |
| /model.10/m/m.0/attn/proj/conv/Conv:         | 28.761%                  |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.0/conv/Conv: | 22.876%                  |
| /model.23/cv2.2/cv2.2.0/conv/Conv:           | 21.570%                  |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.1/conv/Conv: | 21.467%                  |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.1/conv/Conv: | 21.021%                  |
| /model.23/cv2.2/cv2.2.1/conv/Conv:           | 20.973%                  |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.1/conv/Conv: | 19.432%                  |
| /model.22/m.0/cv2/conv/Conv:                 | 19.320%                  |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.1/conv/Conv: | 19.243%                  |
| /model.22/m.0/cv3/conv/Conv:                 | 19.029%                  |
| /model.22/cv2/conv/Conv:                     | 18.488%                  |
| /model.22/m.0/m/m.1/cv2/conv/Conv:           | 18.222%                  |
| /model.23/cv2.1/cv2.1.1/conv/Conv:           | 17.400%                  |
| /model.8/m.0/cv2/conv/Conv:                  | 16.189%                  |
| /model.23/cv2.0/cv2.0.1/conv/Conv:           | 15.585%                  |
| /model.10/m/m.0/attn/pe/conv/Conv:           | 14.687%                  |
| /model.10/m/m.0/attn/qkv/conv/Conv:          | 14.601%                  |
| /model.23/cv2.1/cv2.1.0/conv/Conv:           | 14.154%                  |
| /model.22/cv1/conv/Conv:                     | 14.102%                  |
| /model.10/m/m.0/attn/MatMul_1:               | 13.998%                  |
| /model.10/cv1/conv/Conv:                     | 13.560%                  |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.1/conv/Conv: | 11.771%                  |
| /model.19/m.0/cv2/conv/Conv:                 | 11.216%                  |
| /model.22/m.0/m/m.0/cv2/conv/Conv:           | 11.140%                  |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.0/conv/Conv: | 11.057%                  |
| /model.13/m.0/cv2/conv/Conv:                 | 10.881%                  |
| /model.20/conv/Conv:                         | 10.692%                  |
| /model.23/cv2.2/cv2.2.2/Conv:                | 9.888%                   |
| /model.10/cv2/conv/Conv:                     | 9.788%                   |
| /model.8/cv2/conv/Conv:                      | 9.477%                   |
| /model.8/m.0/cv1/conv/Conv:                  | 9.422%                   |
| /model.19/cv2/conv/Conv:                     | 9.102%                   |
| /model.8/cv1/conv/Conv:                      | 9.101%                   |
| /model.8/m.0/cv3/conv/Conv:                  | 9.068%                   |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.1/conv/Conv: | 9.014%                   |
| /model.22/m.0/m/m.0/cv1/conv/Conv:           | 8.996%                   |
| /model.6/m.0/cv2/conv/Conv:                  | 8.882%                   |
| /model.22/m.0/m/m.1/cv1/conv/Conv:           | 8.637%                   |
| /model.13/cv2/conv/Conv:                     | 8.556%                   |
| /model.8/m.0/m/m.0/cv1/conv/Conv:            | 8.461%                   |
| /model.8/m.0/m/m.0/cv2/conv/Conv:            | 8.362%                   |
| /model.19/cv1/conv/Conv:                     | 8.194%                   |
| /model.8/m.0/m/m.1/cv1/conv/Conv:            | 8.021%                   |
| /model.13/cv1/conv/Conv:                     | 7.910%                   |
| /model.10/m/m.0/attn/MatMul:                 | 7.861%                   |
| /model.19/m.0/cv1/conv/Conv:                 | 7.520%                   |
| /model.22/m.0/cv1/conv/Conv:                 | 7.239%                   |
| /model.8/m.0/m/m.1/cv2/conv/Conv:            | 7.054%                   |
| /model.23/cv2.0/cv2.0.0/conv/Conv:           | 7.042%                   |
| /model.13/m.0/cv1/conv/Conv:                 | 6.987%                   |
| /model.23/cv2.0/cv2.0.2/Conv:                | 6.739%                   |
| /model.23/cv2.1/cv2.1.2/Conv:                | 6.734%                   |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.0/conv/Conv: | 6.660%                   |
| /model.17/conv/Conv:                         | 6.025%                   |
| /model.16/m.0/cv2/conv/Conv:                 | 5.897%                   |
| /model.6/cv2/conv/Conv:                      | 5.815%                   |
| /model.6/m.0/cv3/conv/Conv:                  | 5.814%                   |
| /model.6/cv1/conv/Conv:                      | 5.693%                   |
| /model.7/conv/Conv:                          | 5.570%                   |

(continues on next page)

(continued from previous page)

|  |                        |                           |
|--|------------------------|---------------------------|
| /model.9/cv2/conv/Conv:  | ██                     | 5.382%                    |
| /model.10/m/m.0/ffn/ffn.0/conv/Conv:   | ██                     | 5.173%                    |
| /model.6/m.0/m/m.0/cv1/conv/Conv:  | ██                     | 5.168%                    |
| /model.16/m.0/cv1/conv/Conv:   | ██                     | 5.087%                    |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.0/conv/Conv:   | ██                     | 5.010%                    |
| /model.16/cv2/conv/Conv:   | ██                     | 4.991%                    |
| /model.2/cv2/conv/Conv:  | ██                     | 4.552%                    |
| /model.6/m.0/m/m.0/cv2/conv/Conv:  | ██                     | 4.443%                    |
| /model.3/conv/Conv:  | ██                     | 4.318%                    |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.0/conv/Conv:   | ██                     | 4.304%                    |
| /model.6/m.0/m/m.1/cv1/conv/Conv:  | ██                     | 3.968%                    |
| /model.5/conv/Conv:  | ██                     | 3.948%                    |
| /model.6/m.0/cv1/conv/Conv:  | ██                     | 3.863%                    |
| /model.4/cv1/conv/Conv:  | ██                     | 3.720%                    |
| /model.2/cv1/conv/Conv:  | ██                     | 3.565%                    |
| /model.4/cv2/conv/Conv:  | ██                     | 3.538%                    |
| /model.16/cv1/conv/Conv:   | ██                     | 3.110%                    |
| /model.2/m.0/cv2/conv/Conv:  | ██                     | 2.844%                    |
| /model.6/m.0/m/m.1/cv2/conv/Conv:  | ██                     | 2.762%                    |
| /model.4/m.0/cv1/conv/Conv:  | ██                     | 2.532%                    |
| /model.9/cv1/conv/Conv:  | ██                     | 2.015%                    |
| /model.4/m.0/cv2/conv/Conv:  | ██                     | 1.761%                    |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.0/conv/Conv:   | ██                     | 1.317%                    |
| /model.1/conv/Conv:  | ██                     | 1.315%                    |
| /model.23/cv3.2/cv3.2.2/Conv:  | ██                     | 1.114%                    |
| /model.2/m.0/cv1/conv/Conv:  | ██                     | 0.731%                    |
| /model.23/cv3.1/cv3.1.2/Conv:  | ██                     | 0.491%                    |
| /model.23/cv3.0/cv3.0.2/Conv:  | ██                     | 0.282%                    |
| /model.0/conv/Conv:  | ██                     | 0.159%                    |
| Analysing Layerwise quantization error:: 100% [██████████] 89/89 [07:46<00:00, 5.24s/it] |                        |                           |
| Layer  |                        | NOISE: SIGNAL POWER RATIO |
| /model.1/conv/Conv:  | ██████████████████████ | 0.384%                    |
| /model.22/cv1/conv/Conv:   | ██████████████████     | 0.247%                    |
| /model.4/cv2/conv/Conv:  | ██████████████         | 0.233%                    |
| /model.2/cv2/conv/Conv:  | ██████████             | 0.201%                    |
| /model.0/conv/Conv:  | ██████████             | 0.192%                    |
| /model.9/cv2/conv/Conv:  | ██████████             | 0.156%                    |
| /model.10/cv1/conv/Conv:   | ██████████             | 0.132%                    |
| /model.3/conv/Conv:  | ██████████             | 0.108%                    |
| /model.4/cv1/conv/Conv:  | ██████████             | 0.074%                    |
| /model.16/cv1/conv/Conv:   | ██████████             | 0.066%                    |
| /model.2/cv1/conv/Conv:  | ██████████             | 0.060%                    |
| /model.23/cv2.0/cv2.0.0/conv/Conv:   | ██████████             | 0.052%                    |
| /model.2/m.0/cv1/conv/Conv:  | ██████████             | 0.044%                    |
| /model.6/cv1/conv/Conv:  | ██████████             | 0.033%                    |
| /model.10/m/m.0/attn/pe/conv/Conv:   | ██████████             | 0.029%                    |
| /model.2/m.0/cv2/conv/Conv:  | ██████████             | 0.028%                    |
| /model.22/m.0/m/m.0/cv1/conv/Conv:   | ██████████             | 0.023%                    |
| /model.16/cv2/conv/Conv:   | ██████████             | 0.021%                    |
| /model.16/m.0/cv2/conv/Conv:   | ██████████             | 0.020%                    |
| /model.19/m.0/cv1/conv/Conv:   | ██████████             | 0.020%                    |
| /model.4/m.0/cv1/conv/Conv:  | ██████████             | 0.018%                    |
| /model.19/cv2/conv/Conv:   | ██████████             | 0.017%                    |
| /model.4/m.0/cv2/conv/Conv:  | ██████████             | 0.016%                    |
| /model.10/m/m.0/attn/qkv/conv/Conv:  | ██████████             | 0.016%                    |
| /model.19/cv1/conv/Conv:   | ██████████             | 0.015%                    |
| /model.13/cv2/conv/Conv:   | ██████████             | 0.015%                    |
| /model.8/cv1/conv/Conv:  | ██████████             | 0.013%                    |
| /model.23/cv2.1/cv2.1.0/conv/Conv:   | ██████████             | 0.013%                    |
| /model.23/cv2.2/cv2.2.1/conv/Conv:   | ██████████             | 0.012%                    |
| /model.13/cv1/conv/Conv:   | ██████████             | 0.012%                    |

(continues on next page)

(continued from previous page)

|  |   |        |
|--|---|--------|
| /model.10/cv2/conv/Conv:                     | ■ | 0.011% |
| /model.13/m.0/cv1/conv/Conv:                 | ■ | 0.011% |
| /model.6/cv2/conv/Conv:                      | ■ | 0.011% |
| /model.13/m.0/cv2/conv/Conv:                 | ■ | 0.010% |
| /model.5/conv/Conv:                          |   | 0.010% |
| /model.19/m.0/cv2/conv/Conv:                 |   | 0.009% |
| /model.6/m.0/m.m.1/cv1/conv/Conv:            |   | 0.009% |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.1/conv/Conv: |   | 0.008% |
| /model.23/cv2.2/cv2.2.0/conv/Conv:           |   | 0.008% |
| /model.23/cv2.1/cv2.1.1/conv/Conv:           |   | 0.008% |
| /model.9/cv1/conv/Conv:                      |   | 0.008% |
| /model.23/cv2.0/cv2.0.1/conv/Conv:           |   | 0.007% |
| /model.16/m.0/cv1/conv/Conv:                 |   | 0.007% |
| /model.17/conv/Conv:                         |   | 0.007% |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.0/conv/Conv: |   | 0.007% |
| /model.10/m.m.0/ffn/ffn.1/conv/Conv:         |   | 0.007% |
| /model.23/cv2.0/cv2.0.2/Conv:                |   | 0.006% |
| /model.8/m.0/cv1/conv/Conv:                  |   | 0.006% |
| /model.23/cv2.2/cv2.2.2/Conv:                |   | 0.005% |
| /model.23/cv2.1/cv2.1.2/Conv:                |   | 0.005% |
| /model.22/m.0/cv3/conv/Conv:                 |   | 0.005% |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.1/conv/Conv: |   | 0.005% |
| /model.7/conv/Conv:                          |   | 0.005% |
| /model.8/cv2/conv/Conv:                      |   | 0.004% |
| /model.22/cv2/conv/Conv:                     |   | 0.004% |
| /model.6/m.0/cv3/conv/Conv:                  |   | 0.004% |
| /model.10/m.m.0/ffn/ffn.0/conv/Conv:         |   | 0.004% |
| /model.8/m.0/m.m.1/cv2/conv/Conv:            |   | 0.004% |
| /model.22/m.0/m.m.1/cv1/conv/Conv:           |   | 0.004% |
| /model.8/m.0/m.m.1/cv1/conv/Conv:            |   | 0.004% |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.1/conv/Conv: |   | 0.003% |
| /model.10/m.m.0/attn/proj/conv/Conv:         |   | 0.003% |
| /model.22/m.0/m.m.0/cv2/conv/Conv:           |   | 0.003% |
| /model.22/m.0/cv1/conv/Conv:                 |   | 0.003% |
| /model.8/m.0/cv3/conv/Conv:                  |   | 0.003% |
| /model.6/m.0/m.m.0/cv1/conv/Conv:            |   | 0.003% |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.0/conv/Conv: |   | 0.003% |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.0/conv/Conv: |   | 0.002% |
| /model.6/m.0/m.m.1/cv2/conv/Conv:            |   | 0.002% |
| /model.8/m.0/m.m.0/cv2/conv/Conv:            |   | 0.002% |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.1/conv/Conv: |   | 0.002% |
| /model.10/m.m.0/attn/MatMul_1:               |   | 0.002% |
| /model.22/m.0/m.m.1/cv2/conv/Conv:           |   | 0.001% |
| /model.6/m.0/m.m.0/cv2/conv/Conv:            |   | 0.001% |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.0/conv/Conv: |   | 0.001% |
| /model.8/m.0/m.m.0/cv1/conv/Conv:            |   | 0.001% |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.1/conv/Conv: |   | 0.001% |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.1/conv/Conv: |   | 0.001% |
| /model.6/m.0/cv1/conv/Conv:                  |   | 0.001% |
| /model.23/cv3.2/cv3.2.2/Conv:                |   | 0.001% |
| /model.20/conv/Conv:                         |   | 0.001% |
| /model.23/cv3.1/cv3.1.2/Conv:                |   | 0.001% |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.0/conv/Conv: |   | 0.001% |
| /model.6/m.0/cv2/conv/Conv:                  |   | 0.001% |
| /model.23/cv3.0/cv3.0.2/Conv:                |   | 0.000% |
| /model.10/m.m.0/attn/MatMul:                 |   | 0.000% |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.0/conv/Conv: |   | 0.000% |
| /model.8/m.0/cv2/conv/Conv:                  |   | 0.000% |
| /model.22/m.0/cv2/conv/Conv:                 |   | 0.000% |

**Quantization error analysis**

With the same inputs, The mAP50:95 on COCO val2017 after quantization is only 30.8%, which is lower than that of the float model. There is a accuracy loss with:

- **Graphwise Error**

The output layers of the model are `/model.23/cv3.2/cv3.2.2/Conv`, `/model.23/cv2.2/cv2.2.2/Conv`, `/model.23/cv3.1/cv3.1.2/Conv`, `/model.23/cv2.1/cv2.1.2/Conv`, `/model.23/cv3.0/cv3.0.2/Conv` and `/model.23/cv2.0/cv2.0.2/Conv`. The cumulative error for these layers are 1.114%, 9.888%, 0.491%, 6.734%, 0.282% and 6.739% respectively. Generally, if the cumulative error of the output layer is less than 10%, the loss in accuracy of the quantized model is minimal.

- **Layerwise error**

Observing the Layerwise error, it is found that the errors for all layers are below 1%, indicating that the quantization errors for all layers are small.

We noticed that although the layer-wise errors for all layers are small, the cumulative errors in some layers are relatively large. This may be related to the complex CSP structure in the yolov1n model, where the inputs to the `Concat` or `Add` layers may have different distributions or scales. We can choose to quantize certain layers using `int16` and optimize the quantization with horizontal layer split pass. For more details, please refer to the mixed-precision + horizontal layer split pass quantization test.

## Mixed-Precision + Horizontal Layer Split Pass Quantization

### Quantization settings

```
from ppq.api import get_target_platform
target="esp32p4"
num_of_bits=8
batch_size=32

# Quantize the following layers with 16-bits
quant_setting = QuantizationSettingFactory.espdsl_setting()
quant_setting.dispatching_table.append("/model.2/cv2/conv/Conv", get_target_
↳platform(TARGET, 16))
quant_setting.dispatching_table.append("/model.3/conv/Conv", get_target_
↳platform(TARGET, 16))
quant_setting.dispatching_table.append("/model.4/cv2/conv/Conv", get_target_
↳platform(TARGET, 16))

# Horizontal Layer Split Pass
quant_setting.weight_split = True
quant_setting.weight_split_setting.method = 'balance'
quant_setting.weight_split_setting.value_threshold = 1.5
quant_setting.weight_split_setting.interested_layers = ['/model.0/conv/Conv', '/
↳model.1/conv/Conv']
```

### Quantization results

| Layer  | NOISE:SIGNAL POWER RATIO |
|--|--------------------------|
| <code>/model.10/m/m.0/ffn/ffn.1/conv/Conv</code> :         | 24.841%                  |
| <code>/model.10/m/m.0/attn/proj/conv/Conv</code> :         | 19.061%                  |
| <code>/model.23/cv2.2/cv2.2.1/conv/Conv</code> :           | 17.927%                  |
| <code>/model.23/cv3.2/cv3.2.0/cv3.2.0.0/conv/Conv</code> : | 17.396%                  |
| <code>/model.23/cv2.2/cv2.2.0/conv/Conv</code> :           | 17.061%                  |
| <code>/model.22/m.0/cv3/conv/Conv</code> :                 | 15.563%                  |
| <code>/model.23/cv3.2/cv3.2.0/cv3.2.0.1/conv/Conv</code> : | 15.427%                  |
| <code>/model.23/cv3.0/cv3.0.1/cv3.0.1.1/conv/Conv</code> : | 14.890%                  |
| <code>/model.22/m.0/m/m.1/cv2/conv/Conv</code> :           | 14.784%                  |
| <code>/model.23/cv3.2/cv3.2.1/cv3.2.1.1/conv/Conv</code> : | 14.243%                  |
| <code>/model.22/cv2/conv/Conv</code> :                     | 14.098%                  |
| <code>/model.22/m.0/cv2/conv/Conv</code> :                 | 13.945%                  |
| <code>/model.23/cv3.1/cv3.1.1/cv3.1.1.1/conv/Conv</code> : | 13.489%                  |
| <code>/model.23/cv2.1/cv2.1.1/conv/Conv</code> :           | 10.919%                  |

(continues on next page)



(continued from previous page)

|  |            |         |
|--|------------|---------|
| /model.23/cv2.0/cv2.0.1/conv/Conv:           | ██████████ | 10.073% |
| /model.23/cv2.1/cv2.1.0/conv/Conv:           | ██████████ | 9.819%  |
| /model.22/cv1/conv/Conv:                     | ██████████ | 9.093%  |
| /model.10/m/m.0/attn/MatMul_1:               | ██████████ | 8.414%  |
| /model.22/m.0/m/m.0/cv2/conv/Conv:           | ██████████ | 8.245%  |
| /model.23/cv2.2/cv2.2.2/Conv:                | ██████████ | 8.208%  |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.1/conv/Conv: | ██████████ | 8.031%  |
| /model.10/m/m.0/attn/qkv/conv/Conv:          | ██████████ | 7.818%  |
| /model.13/m.0/cv2/conv/Conv:                 | ██████████ | 7.717%  |
| /model.19/m.0/cv2/conv/Conv:                 | ██████████ | 7.404%  |
| /model.20/conv/Conv:                         | ██████████ | 7.161%  |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.0/conv/Conv: | ██████████ | 7.080%  |
| /model.10/m/m.0/attn/pe/conv/Conv:           | ██████████ | 6.814%  |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.1/conv/Conv: | ██████████ | 6.764%  |
| /model.22/m.0/m/m.1/cv1/conv/Conv:           | ██████████ | 6.539%  |
| /model.22/m.0/m/m.0/cv1/conv/Conv:           | ██████████ | 6.418%  |
| /model.19/cv2/conv/Conv:                     | ██████████ | 6.206%  |
| /model.13/cv2/conv/Conv:                     | ██████████ | 5.894%  |
| /model.10/cv1/conv/Conv:                     | ██████████ | 5.757%  |
| /model.10/cv2/conv/Conv:                     | ██████████ | 5.716%  |
| /model.19/cv1/conv/Conv:                     | ██████████ | 5.279%  |
| /model.22/m.0/cv1/conv/Conv:                 | ██████████ | 5.072%  |
| /model.19/m.0/cv1/conv/Conv:                 | ██████████ | 5.036%  |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.0/conv/Conv: | ██████████ | 4.979%  |
| /model.8/m.0/cv2/conv/Conv:                  | ██████████ | 4.862%  |
| /model.10/m/m.0/attn/MatMul:                 | ██████████ | 4.670%  |
| /model.13/cv1/conv/Conv:                     | ██████████ | 4.594%  |
| /model.23/cv2.0/cv2.0.0/conv/Conv:           | ██████████ | 4.441%  |
| /model.23/cv2.0/cv2.0.2/Conv:                | ██████████ | 4.308%  |
| /model.13/m.0/cv1/conv/Conv:                 | ██████████ | 4.278%  |
| /model.23/cv2.1/cv2.1.2/Conv:                | ██████████ | 4.214%  |
| /model.6/m.0/cv2/conv/Conv:                  | ██████████ | 4.031%  |
| /model.17/conv/Conv:                         | ██████████ | 3.760%  |
| /model.16/m.0/cv2/conv/Conv:                 | ██████████ | 3.521%  |
| /model.8/m.0/cv1/conv/Conv:                  | ██████████ | 3.227%  |
| /model.16/m.0/cv1/conv/Conv:                 | ██████████ | 3.185%  |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.0/conv/Conv: | ██████████ | 3.178%  |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.0/conv/Conv: | ██████████ | 3.150%  |
| /model.8/cv2/conv/Conv:                      | ██████████ | 3.067%  |
| /model.8/m.0/cv3/conv/Conv:                  | ██████████ | 3.067%  |
| /model.16/cv2/conv/Conv:                     | ██████████ | 3.054%  |
| /model.2/cv2/conv/Conv:                      | ██████████ | 3.053%  |
| /model.8/m.0/m/m.1/cv1/conv/Conv:            | ██████████ | 3.049%  |
| /model.6/m.0/cv3/conv/Conv:                  | ██████████ | 3.049%  |
| /model.8/cv1/conv/Conv:                      | ██████████ | 2.984%  |
| /model.8/m.0/m/m.0/cv2/conv/Conv:            | ██████████ | 2.934%  |
| /model.10/m/m.0/ffn/ffn.0/conv/Conv:         | ██████████ | 2.794%  |
| /model.6/cv1/conv/Conv:                      | ██████████ | 2.783%  |
| /model.8/m.0/m/m.0/cv1/conv/Conv:            | ██████████ | 2.753%  |
| /model.2/cv1/conv/Conv:                      | ██████████ | 2.697%  |
| /model.6/cv2/conv/Conv:                      | ██████████ | 2.616%  |
| /model.8/m.0/m/m.1/cv2/conv/Conv:            | ██████████ | 2.596%  |
| /model.9/cv2/conv/Conv:                      | ██████████ | 2.500%  |
| /model.3/conv/Conv:                          | ██████████ | 2.499%  |
| /model.2/m.0/cv2/conv/Conv:                  | ██████████ | 2.469%  |
| /model.6/m.0/m/m.0/cv2/conv/Conv:            | ██████████ | 2.235%  |
| /model.6/m.0/m/m.0/cv1/conv/Conv:            | ██████████ | 2.233%  |
| /model.4/cv2/conv/Conv:                      | ██████████ | 2.150%  |
| /model.7/conv/Conv:                          | ██████████ | 2.075%  |
| /model.6/m.0/m/m.1/cv1/conv/Conv:            | ██████████ | 2.069%  |
| /model.5/conv/Conv:                          | ██████████ | 1.998%  |

(continues on next page)



(continued from previous page)

|  |  |        |
|--|--|--------|
| /model.10/m/m.0/ffn/ffn.1/conv/Conv:         |  | 0.007% |
| /model.22/m.0/cv1/conv/Conv:                 |  | 0.006% |
| /model.10/cv2/conv/Conv:                     |  | 0.006% |
| /model.23/cv2.0/cv2.0.2/Conv:                |  | 0.006% |
| /model.23/cv2.2/cv2.2.2/Conv:                |  | 0.005% |
| /model.23/cv2.1/cv2.1.2/Conv:                |  | 0.005% |
| /model.22/m.0/cv3/conv/Conv:                 |  | 0.005% |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.1/conv/Conv: |  | 0.005% |
| /model.22/cv2/conv/Conv:                     |  | 0.005% |
| /model.7/conv/Conv:                          |  | 0.004% |
| /model.6/m.0/cv3/conv/Conv:                  |  | 0.004% |
| /model.10/m/m.0/ffn/ffn.0/conv/Conv:         |  | 0.004% |
| /model.8/m.0/m/m.1/cv2/conv/Conv:            |  | 0.004% |
| /model.22/m.0/m/m.1/cv1/conv/Conv:           |  | 0.004% |
| /model.8/m.0/m/m.1/cv1/conv/Conv:            |  | 0.004% |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.1/conv/Conv: |  | 0.003% |
| /model.8/m.0/cv1/conv/Conv:                  |  | 0.003% |
| /model.10/m/m.0/attn/proj/conv/Conv:         |  | 0.003% |
| /model.22/m.0/m/m.0/cv2/conv/Conv:           |  | 0.003% |
| PPQ_Operation_2:                             |  | 0.003% |
| /model.8/m.0/cv3/conv/Conv:                  |  | 0.003% |
| /model.6/m.0/m/m.0/cv1/conv/Conv:            |  | 0.003% |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.0/conv/Conv: |  | 0.002% |
| /model.6/m.0/m/m.1/cv2/conv/Conv:            |  | 0.002% |
| /model.8/m.0/m/m.0/cv2/conv/Conv:            |  | 0.002% |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.0/conv/Conv: |  | 0.002% |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.1/conv/Conv: |  | 0.002% |
| /model.10/m/m.0/attn/MatMul_1:               |  | 0.002% |
| /model.22/m.0/m/m.1/cv2/conv/Conv:           |  | 0.001% |
| /model.6/m.0/m/m.0/cv2/conv/Conv:            |  | 0.001% |
| /model.8/m.0/m/m.0/cv1/conv/Conv:            |  | 0.001% |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.0/conv/Conv: |  | 0.001% |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.1/conv/Conv: |  | 0.001% |
| /model.2/cv2/conv/Conv:                      |  | 0.001% |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.1/conv/Conv: |  | 0.001% |
| /model.6/m.0/cv1/conv/Conv:                  |  | 0.001% |
| /model.23/cv3.2/cv3.2.2/Conv:                |  | 0.001% |
| /model.20/conv/Conv:                         |  | 0.001% |
| /model.23/cv3.1/cv3.1.2/Conv:                |  | 0.001% |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.0/conv/Conv: |  | 0.001% |
| /model.6/m.0/cv2/conv/Conv:                  |  | 0.001% |
| /model.23/cv3.0/cv3.0.2/Conv:                |  | 0.000% |
| /model.10/m/m.0/attn/MatMul:                 |  | 0.000% |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.0/conv/Conv: |  | 0.000% |
| /model.8/m.0/cv2/conv/Conv:                  |  | 0.000% |
| /model.22/m.0/cv2/conv/Conv:                 |  | 0.000% |
| /model.3/conv/Conv:                          |  | 0.000% |
| /model.4/cv2/conv/Conv:                      |  | 0.000% |

### Quantization error analysis

After using 16-bits quantization on layers with higher layer-wise error and employing horizontal layer split pass, the quantized model's mAP50:95 on COCO val2017 improves to 33.4% with the same inputs. Additionally, a noticeable decrease in cumulative error of output layers can be observed.

The graphwise error for the output layers of the model, /model.23/cv3.2/cv3.2.2/Conv, /model.23/cv2.2/cv2.2.2/Conv, /model.23/cv3.1/cv3.1.2/Conv, /model.23/cv2.1/cv2.1.2/Conv, /model.23/cv3.0/cv3.0.2/Conv and /model.23/cv2.0/cv2.0.2/Conv, are 0.784%, 8.208%, 0.360%, 4.214%, 0.189% and 4.308% respectively.

## Quantization-Aware Training

To further improve the accuracy of the quantized model, we adopt the quantization-aware training(QAT) strategy. Here, QAT is performed based on 8-bit quantization.

### Quantization settings

- [yolo11n\\_qat.py](#)
- [train.py](#)

### Quantization results

| Layer  | NOISE:SIGNAL POWER RATIO |
|--|--------------------------|
| /model.10/m/m.0/ffn/ffn.1/conv/Conv:         | 23.754%                  |
| /model.10/m/m.0/attn/proj/conv/Conv:         | 16.118%                  |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.1/conv/Conv: | 10.878%                  |
| /model.8/m.0/cv2/conv/Conv:                  | 10.527%                  |
| /model.22/m.0/cv3/conv/Conv:                 | 10.298%                  |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.1/conv/Conv: | 10.188%                  |
| /model.10/m/m.0/attn/pe/conv/Conv:           | 10.093%                  |
| /model.22/m.0/m/m.1/cv2/conv/Conv:           | 9.891%                   |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.0/conv/Conv: | 9.839%                   |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.1/conv/Conv: | 9.827%                   |
| /model.23/cv2.2/cv2.2.0/conv/Conv:           | 9.658%                   |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.1/conv/Conv: | 9.168%                   |
| /model.22/m.0/cv2/conv/Conv:                 | 8.604%                   |
| /model.10/m/m.0/attn/MatMul_1:               | 8.596%                   |
| /model.10/m/m.0/attn/qkv/conv/Conv:          | 8.541%                   |
| /model.23/cv2.2/cv2.2.1/conv/Conv:           | 8.528%                   |
| /model.22/cv2/conv/Conv:                     | 8.442%                   |
| /model.23/cv2.1/cv2.1.1/conv/Conv:           | 8.306%                   |
| /model.23/cv2.0/cv2.0.1/conv/Conv:           | 8.015%                   |
| /model.10/cv1/conv/Conv:                     | 7.998%                   |
| /model.22/cv1/conv/Conv:                     | 7.307%                   |
| /model.8/cv1/conv/Conv:                      | 7.265%                   |
| /model.23/cv2.1/cv2.1.0/conv/Conv:           | 6.989%                   |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.1/conv/Conv: | 6.716%                   |
| /model.6/m.0/cv2/conv/Conv:                  | 6.595%                   |
| /model.2/cv2/conv/Conv:                      | 6.131%                   |
| /model.22/m.0/m/m.0/cv2/conv/Conv:           | 6.078%                   |
| /model.10/m/m.0/attn/MatMul:                 | 6.055%                   |
| /model.19/m.0/cv2/conv/Conv:                 | 5.999%                   |
| /model.8/m.0/cv1/conv/Conv:                  | 5.919%                   |
| /model.13/m.0/cv2/conv/Conv:                 | 5.863%                   |
| /model.20/conv/Conv:                         | 5.638%                   |
| /model.8/cv2/conv/Conv:                      | 5.616%                   |
| /model.10/cv2/conv/Conv:                     | 5.464%                   |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.1/conv/Conv: | 5.443%                   |
| /model.2/m.0/cv2/conv/Conv:                  | 5.426%                   |
| /model.8/m.0/m/m.0/cv1/conv/Conv:            | 5.390%                   |
| /model.13/cv2/conv/Conv:                     | 5.256%                   |
| /model.19/cv2/conv/Conv:                     | 5.231%                   |
| /model.13/cv1/conv/Conv:                     | 5.131%                   |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.0/conv/Conv: | 5.122%                   |
| /model.6/cv1/conv/Conv:                      | 5.049%                   |
| /model.6/cv2/conv/Conv:                      | 4.788%                   |
| /model.8/m.0/m/m.0/cv2/conv/Conv:            | 4.706%                   |
| /model.19/cv1/conv/Conv:                     | 4.586%                   |
| /model.7/conv/Conv:                          | 4.586%                   |
| /model.8/m.0/m/m.1/cv1/conv/Conv:            | 4.541%                   |
| /model.8/m.0/cv3/conv/Conv:                  | 4.529%                   |
| /model.3/conv/Conv:                          | 4.361%                   |
| /model.13/m.0/cv1/conv/Conv:                 | 4.359%                   |

(continues on next page)

(continued from previous page)

|   |                          |        |
|---|--------------------------|--------|
| /model.22/m.0/m/m.1/cv1/conv/Conv:  |                          | 4.328% |
| /model.6/m.0/cv3/conv/Conv:   |                          | 4.156% |
| /model.22/m.0/m/m.0/cv1/conv/Conv:  |                          | 4.083% |
| /model.23/cv2.0/cv2.0.0/conv/Conv:  |                          | 3.998% |
| /model.19/m.0/cv1/conv/Conv:  |                          | 3.974% |
| /model.23/cv2.2/cv2.2.2/Conv:   |                          | 3.817% |
| /model.16/m.0/cv1/conv/Conv:  |                          | 3.797% |
| /model.16/m.0/cv2/conv/Conv:  |                          | 3.654% |
| /model.4/cv1/conv/Conv:   |                          | 3.544% |
| /model.4/cv2/conv/Conv:   |                          | 3.488% |
| /model.22/m.0/cv1/conv/Conv:  |                          | 3.423% |
| /model.8/m.0/m/m.1/cv2/conv/Conv:   |                          | 3.382% |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.0/conv/Conv:                                    |                          | 3.299% |
| /model.17/conv/Conv:  |                          | 3.296% |
| /model.6/m.0/m/m.0/cv1/conv/Conv:   |                          | 3.267% |
| /model.5/conv/Conv:   |                          | 3.147% |
| /model.23/cv2.1/cv2.1.2/Conv:   |                          | 3.102% |
| /model.16/cv2/conv/Conv:  |                          | 3.091% |
| /model.6/m.0/m/m.0/cv2/conv/Conv:   |                          | 3.080% |
| /model.23/cv2.0/cv2.0.2/Conv:   |                          | 3.056% |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.0/conv/Conv:                                    |                          | 2.989% |
| /model.2/cv1/conv/Conv:   |                          | 2.874% |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.0/conv/Conv:                                    |                          | 2.843% |
| /model.6/m.0/cv1/conv/Conv:   |                          | 2.819% |
| /model.9/cv2/conv/Conv:   |                          | 2.662% |
| /model.6/m.0/m/m.1/cv1/conv/Conv:   |                          | 2.633% |
| /model.10/m/m.0/ffn/ffn.0/conv/Conv:  |                          | 2.581% |
| /model.4/m.0/cv1/conv/Conv:   |                          | 2.545% |
| /model.16/cv1/conv/Conv:  |                          | 2.171% |
| /model.4/m.0/cv2/conv/Conv:   |                          | 1.942% |
| /model.6/m.0/m/m.1/cv2/conv/Conv:   |                          | 1.925% |
| /model.2/m.0/cv1/conv/Conv:   |                          | 1.721% |
| /model.9/cv1/conv/Conv:   |                          | 1.140% |
| /model.1/conv/Conv:   |                          | 1.117% |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.0/conv/Conv:                                    |                          | 0.831% |
| /model.23/cv3.2/cv3.2.2/Conv:   |                          | 0.443% |
| /model.23/cv3.1/cv3.1.2/Conv:   |                          | 0.247% |
| /model.0/conv/Conv:   |                          | 0.150% |
| /model.23/cv3.0/cv3.0.2/Conv:   |                          | 0.119% |
| Analysing Layerwise quantization error:: 100%    89/89 [04:44<00:00, 3.↵20s/it] |                          |        |
| Layer   | NOISE:SIGNAL POWER RATIO |        |
| /model.2/cv2/conv/Conv:   |                          | 1.462% |
| /model.3/conv/Conv:   |                          | 0.764% |
| /model.4/cv2/conv/Conv:   |                          | 0.763% |
| /model.10/cv2/conv/Conv:  |                          | 0.535% |
| /model.9/cv2/conv/Conv:   |                          | 0.439% |
| /model.2/cv1/conv/Conv:   |                          | 0.395% |
| /model.4/cv1/conv/Conv:   |                          | 0.361% |
| /model.1/conv/Conv:   |                          | 0.347% |
| /model.2/m.0/cv1/conv/Conv:   |                          | 0.192% |
| /model.4/m.0/cv2/conv/Conv:   |                          | 0.184% |
| /model.22/cv1/conv/Conv:  |                          | 0.179% |
| /model.5/conv/Conv:   |                          | 0.161% |
| /model.16/cv1/conv/Conv:  |                          | 0.154% |
| /model.10/cv1/conv/Conv:  |                          | 0.145% |
| /model.16/m.0/cv2/conv/Conv:  |                          | 0.142% |
| /model.16/m.0/cv1/conv/Conv:  |                          | 0.113% |
| /model.4/m.0/cv1/conv/Conv:   |                          | 0.107% |
| /model.0/conv/Conv:   |                          | 0.100% |
| /model.10/m/m.0/attn/pe/conv/Conv:  |                          | 0.095% |

(continues on next page)

(continued from previous page)

|  |   |        |
|--|---|--------|
| /model.6/cv1/conv/Conv:                      | █ | 0.082% |
| /model.23/cv2.2/cv2.2.2/Conv:                | █ | 0.082% |
| /model.16/cv2/conv/Conv:                     | █ | 0.076% |
| /model.6/cv2/conv/Conv:                      | █ | 0.066% |
| /model.22/m.0/cv1/conv/Conv:                 | █ | 0.060% |
| /model.13/cv2/conv/Conv:                     | █ | 0.056% |
| /model.19/cv2/conv/Conv:                     | █ | 0.041% |
| /model.10/m/m.0/attn/qkv/conv/Conv:          |   | 0.034% |
| /model.7/conv/Conv:                          |   | 0.033% |
| /model.13/cv1/conv/Conv:                     |   | 0.033% |
| /model.23/cv2.2/cv2.2.0/conv/Conv:           |   | 0.032% |
| /model.10/m/m.0/ffn/ffn.0/conv/Conv:         |   | 0.032% |
| /model.23/cv2.0/cv2.0.0/conv/Conv:           |   | 0.029% |
| /model.13/m.0/cv1/conv/Conv:                 |   | 0.029% |
| /model.2/m.0/cv2/conv/Conv:                  |   | 0.026% |
| /model.19/cv1/conv/Conv:                     |   | 0.025% |
| /model.6/m.0/cv3/conv/Conv:                  |   | 0.024% |
| /model.19/m.0/cv2/conv/Conv:                 |   | 0.024% |
| /model.17/conv/Conv:                         |   | 0.023% |
| /model.23/cv2.0/cv2.0.2/Conv:                |   | 0.021% |
| /model.19/m.0/cv1/conv/Conv:                 |   | 0.019% |
| /model.23/cv3.2/cv3.2.2/Conv:                |   | 0.019% |
| /model.9/cv1/conv/Conv:                      |   | 0.017% |
| /model.23/cv2.1/cv2.1.0/conv/Conv:           |   | 0.015% |
| /model.8/cv1/conv/Conv:                      |   | 0.014% |
| /model.22/m.0/cv3/conv/Conv:                 |   | 0.014% |
| /model.13/m.0/cv2/conv/Conv:                 |   | 0.014% |
| /model.8/m.0/cv3/conv/Conv:                  |   | 0.012% |
| /model.23/cv2.2/cv2.2.1/conv/Conv:           |   | 0.011% |
| /model.23/cv2.1/cv2.1.2/Conv:                |   | 0.011% |
| /model.22/m.0/m/m.1/cv1/conv/Conv:           |   | 0.010% |
| /model.22/m.0/m/m.0/cv1/conv/Conv:           |   | 0.009% |
| /model.20/conv/Conv:                         |   | 0.009% |
| /model.8/cv2/conv/Conv:                      |   | 0.009% |
| /model.6/m.0/m/m.1/cv1/conv/Conv:            |   | 0.008% |
| /model.10/m/m.0/ffn/ffn.1/conv/Conv:         |   | 0.008% |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.1/conv/Conv: |   | 0.008% |
| /model.23/cv2.1/cv2.1.1/conv/Conv:           |   | 0.008% |
| /model.23/cv2.0/cv2.0.1/conv/Conv:           |   | 0.007% |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.1/conv/Conv: |   | 0.007% |
| /model.10/m/m.0/attn/proj/conv/Conv:         |   | 0.007% |
| /model.8/m.0/m/m.1/cv1/conv/Conv:            |   | 0.007% |
| /model.8/m.0/cv1/conv/Conv:                  |   | 0.007% |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.0/conv/Conv: |   | 0.006% |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.1/conv/Conv: |   | 0.005% |
| /model.22/cv2/conv/Conv:                     |   | 0.005% |
| /model.6/m.0/m/m.0/cv1/conv/Conv:            |   | 0.004% |
| /model.22/m.0/m/m.0/cv2/conv/Conv:           |   | 0.004% |
| /model.23/cv3.1/cv3.1.1/cv3.1.1.1/conv/Conv: |   | 0.003% |
| /model.6/m.0/cv1/conv/Conv:                  |   | 0.003% |
| /model.8/m.0/m/m.0/cv1/conv/Conv:            |   | 0.003% |
| /model.8/m.0/m/m.1/cv2/conv/Conv:            |   | 0.003% |
| /model.8/m.0/m/m.0/cv2/conv/Conv:            |   | 0.003% |
| /model.6/m.0/m/m.1/cv2/conv/Conv:            |   | 0.003% |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.0/conv/Conv: |   | 0.002% |
| /model.23/cv3.1/cv3.1.2/Conv:                |   | 0.002% |
| /model.23/cv3.0/cv3.0.0/cv3.0.0.0/conv/Conv: |   | 0.002% |
| /model.23/cv3.2/cv3.2.1/cv3.2.1.1/conv/Conv: |   | 0.002% |
| /model.22/m.0/m/m.1/cv2/conv/Conv:           |   | 0.002% |
| /model.6/m.0/m/m.0/cv2/conv/Conv:            |   | 0.002% |
| /model.10/m/m.0/attn/MatMul_1:               |   | 0.002% |

(continues on next page)

(continued from previous page)

|  |  |        |
|--|--|--------|
| /model.23/cv3.0/cv3.0.2/Conv:                |  | 0.001% |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.0/conv/Conv: |  | 0.001% |
| /model.23/cv3.0/cv3.0.1/cv3.0.1.1/conv/Conv: |  | 0.001% |
| /model.23/cv3.2/cv3.2.0/cv3.2.0.0/conv/Conv: |  | 0.001% |
| /model.6/m.0/cv2/conv/Conv:                  |  | 0.000% |
| /model.23/cv3.1/cv3.1.0/cv3.1.0.0/conv/Conv: |  | 0.000% |
| /model.10/m/m.0/attn/MatMul:                 |  | 0.000% |
| /model.8/m.0/cv2/conv/Conv:                  |  | 0.000% |
| /model.22/m.0/cv2/conv/Conv:                 |  | 0.000% |

### Quantization error analysis

After applying QAT to 8-bit quantization, the quantized model's mAP50:95 on COCO val2017 improves to 35.5% with the same inputs, while cumulative errors of out layers are significantly reduced. Compared to the other two quantization methods, the 8-bit QAT quantized model achieves the highest quantization accuracy with the lowest inference latency.

The graphwise error for the output layers of the model, /model.23/cv3.2/cv3.2.2/Conv, /model.23/cv2.2/cv2.2.2/Conv, /model.23/cv3.1/cv3.1.2/Conv, /model.23/cv2.1/cv2.1.2/Conv, /model.23/cv3.0/cv3.0.2/Conv and /model.23/cv2.0/cv2.0.2/Conv, are 0.443%, 3.817%, 0.247%, 3.102%, 0.119% and 3.056% respectively.

### 3.5.3 Model deployment

example

#### Object detection base class

- [dl\\_detect\\_base.hpp](#)
- [dl\\_detect\\_base.cpp](#)

#### Pre-process

ImagePreprocessor class contains the common pre-precoess pipeline, color conversion, crop, re-size, normalization, quantize.

- [dl\\_image\\_preprocessor.hpp](#)
- [dl\\_image\\_preprocessor.cpp](#)

#### Post-process

- [dl\\_detect\\_postprocessor.hpp](#)
- [dl\\_detect\\_postprocessor.cpp](#)
- [dl\\_detect\\_yolo11\\_postprocessor.hpp](#)
- [dl\\_detect\\_yolo11\\_postprocessor.cpp](#)

## 3.6 Creating a New Module (Operator)

This tutorial guides you through the process of creating a new module in the `dl::module` namespace. The `Module` class serves as the base class for all modules, and you can extend this base class to create your custom module.

---

**Note:** The interface of modules in ESP-DL should be aligned with ONNX.

---

### 3.6.1 Understand the Base Module Class

The base class provides several virtual methods that must be overridden in your derived class.

- **Methods:**

- `dl::module::Module::Module()`: Constructor to initialize the module.
- `dl::module::Module::~~Module()`: Destructor to release resources.
- `dl::module::Module::get_output_shape()`: Calculates the output shape based on the input shape.
- `dl::module::Module::forward()`: Runs the module, high-level interface.
- `dl::module::Module::forward_args()`: Runs the module, low-level interface.
- `dl::module::Module::deserialize()`: Creates a module instance from serialized information.
- `dl::module::Module::print()`: Prints module information.

For more information, please refer to [Module Class Reference](#).

### 3.6.2 Create a New Module Class

To create a new module, you need to derive a new class from the `Module` base class and override the necessary methods.

#### Example: Creating a `MyCustomModule` Class

For more examples, please refer to [esp-dl/dl/module](#).

```
#include "module.h" // Include the header file where the Module class is defined

namespace dl {
namespace module {

class MyCustomModule : public Module {
public:
    // Constructor
    MyCustomModule(const char *name = "MyCustomModule",
                  module_inplace_t inplace = MODULE_NON_INPLACE,
                  quant_type_t quant_type = QUANT_TYPE_NONE)
        : Module(name, inplace, quant_type) {}

    // Destructor
    virtual ~MyCustomModule() {}

    // Override the get_output_shape method
    std::vector<std::vector<int>> get_output_shape(std::vector<std::vector<int>> &
↪input_shapes) override {
        // Implement the logic to calculate the output shape based on input shapes
        std::vector<std::vector<int>> output_shapes;
        // Example: Assume the output shape is the same as the input shape
        output_shapes.push_back(input_shapes[0]);
        return output_shapes;
    }

    // Override the forward method
    void forward(std::vector<dl::TensorBase * > &tensors, runtime_mode_t mode =_
↪RUNTIME_MODE_AUTO) override {
```

(continues on next page)



```

    // Implement the logic to run the module
    // Example: Perform some operation on the tensors
    for (auto &tensor : tensors) {
        // Perform some operation on each tensor
    }
}

// Override the forward_args method
void forward_args(void *args) override {
    // Implement the low-level interface logic
    // Example: Perform some operation based on the arguments
}

// Deserialize module instance by serialization information
static Module *deserialize(fbs::FbsModel *fbs_model, std::string node_name){
    // Implement the logic to deserialize the module instance
    // The interface should be align with ONNX
}

// Override the print method
void print() override {
    // Print module information
    ESP_LOGI("MyCustomModule", "Module Name: %s, Quant type: %d", name.c_str(),
↪ quant_type);
}
};

} // namespace module
} // namespace dl

```

### Register MyCustomModule Class

Once you have implemented MyCustomModule Class, register your module in `dl_module_creator` as a globally available module.

```

void register_dl_modules()
{
    if (creators.empty()) {
        ...
        this->register_module("MyCustomModule", MyCustomModule::deserialize);
    }
}

```

# Chapter 4

## API Reference

### 4.1 Tensor API Reference

Tensor is the fundamental data type in esp-dl, used for storing multi-type data such as int8, int16, float, etc., similar to the tensor in PyTorch. We have implemented some common tensor operations. Please refer to the following APIs for details.

#### 4.1.1 Header File

- [esp-dl/dl/tensor/include/dl\\_tensor\\_base.hpp](#)

#### 4.1.2 Classes

class **TensorBase**

This class is designed according to PyTorch Tensor. *TensorBase* is required to ensure that the first address are aligned to 16 bytes and the memory size should be a multiple of 16 bytes.

TODO:: Implement more functions

#### Public Functions

**TensorBase** (std::vector<int> shape, const void \*element, int exponent = 0, dtype\_t dtype = DATA\_TYPE\_FLOAT, bool deep = true, uint32\_t caps = MALLOC\_CAP\_DEFAULT)

Construct a *TensorBase* object.

#### Parameters

- **shape** –Shape of tensor
- **element** –Pointer of data
- **exponent** –Exponent of tensor, default is 0
- **dtype** –Data type of element, default is float
- **deep** –True: malloc memory and copy data, false: use the pointer directly
- **caps** –Bitwise OR of MALLOC\_CAP\_\* flags indicating the type of memory to be returned

inline virtual ~**TensorBase** ()

Destroy the *TensorBase* object.

bool **assign** (*TensorBase* \*tensor)

Assign tensor to this tensor.

**Parameters** **tensor** –

**Returns** true if assign successfully, otherwise false.

bool **assign** (std::vector<int> shape, const void \*element, int exponent, dtype\_t dtype)

Assign data to this tensor.

**Parameters**

- **shape** –
- **element** –
- **exponent** –
- **dtype** –

**Returns** true if assign successfully, otherwise false.

inline int **get\_size** ()

Get the size of Tensor.

**Returns** the size of Tensor.

inline int **get\_aligned\_size** ()

Get the aligned size of Tensor.

**Returns** the aligned size of Tensor.

inline size\_t **get\_dtype\_bytes** ()

Get the dtype size, in bytes.

**Returns** the size of dtype.

inline const char \***get\_dtype\_string** ()

Get the dtype string of Tensor.

**Returns** the string of Tensor's dtype.

inline int **get\_bytes** ()

Get the bytes of Tensor.

**Returns** the bytes of Tensor.

inline int **get\_aligned\_bytes** ()

Get the bytes of Tensor.

**Returns** the bytes of Tensor.

inline virtual void \***get\_element\_ptr** ()

Get data pointer. If cache(preload data pointer) is not null, return cache pointer, otherwise return data pointer.

**Returns** the pointer of Tensor's data

template<typename T>

inline T \***get\_element\_ptr** ()

Get data pointer by the specified template. If cache(preload data pointer) is not null, return cache pointer, otherwise return data pointer.

**Returns** the pointer of Tensor's data

*TensorBase* &**set\_element\_ptr** (void \*data)

Set the data pointer of Tensor.

**Parameters** **data** –point to data memory

**Returns** *TensorBase*& self

```
inline std::vector<int> get_shape ()
```

Get the shape of Tensor.

**Returns** std::vector<int> the shape of Tensor

```
TensorBase &set_shape (const std::vector<int> shape)
```

Set the shape of Tensor.

**Parameters** **shape** –the shape of Tensor.

**Returns** Tensor.

```
inline int get_exponent ()
```

Get the exponent of Tensor.

**Returns** int the exponent of Tensor

```
inline dtype_t get_dtype ()
```

Get the data type of Tensor.

**Returns** dtype\_t the data type of Tensor

```
inline uint32_t get_caps ()
```

Get the memory flags of Tensor.

**Returns** uint32\_t the memory flags of Tensor

```
TensorBase &reshape (std::vector<int> shape)
```

Change a new shape to the Tensor without changing its data.

**Parameters** **shape** –the target shape

**Returns** *TensorBase*& self

```
template<typename T>
```

```
TensorBase *flip (const std::vector<int> &axes)
```

Flip the input Tensor along the specified axes.

**Parameters** **axes** –the specified axes

**Returns** *TensorBase*& self

```
TensorBase *transpose (TensorBase *input, std::vector<int> perm = {})
```

Reverse or permute the axes of the input Tensor.

**Parameters**

- **input** –the input Tensor
- **perm** –the new arrangement of the dims. if perm == {}, the dims arrangement will be reversed.

**Returns** *TensorBase* \*self

```
template<typename T>
```

```
TensorBase *transpose (T *input_element, std::vector<int> &input_shape, std::vector<int> &input_axis_offset, std::vector<int> &perm)
```

Reverse or permute the axes of the input Tensor.

**Parameters**

- **input\_element** –the input data pointer
- **input\_shape** –the input data shape
- **input\_axis\_offset** –the input data axis offset
- **perm** –the new arrangement of the dims. if perm == {}, the dims arrangement will be reversed.

**Returns** *TensorBase* \*self

```
bool is_same_shape (TensorBase *tensor)
```

Check the shape is the same as the shape of input.

**Parameters** **tensor** –Input tensor pointer

**Returns**

- true: same shape
- false: not

bool **equal** (*TensorBase* \*tensor, float epsilon = 1e-6, bool verbose = false)

Compare the shape and data of two Tensor.

**Parameters**

- **tensor** –Input tensor
- **epsilon** –The max error of two element
- **verbose** –If true, print the detail of results

**Returns** true if two tensor is equal otherwise false

*TensorBase* \***slice** (const std::vector<int> &start, const std::vector<int> &end, const std::vector<int> &axes = {}, const std::vector<int> &step = {})

Produces a slice of the this tensor along multiple axes.

**Warning:** The length of start, end and step must be same as the shape of input tensor

**Parameters**

- **start** –Starting indices
- **end** –Ending indices
- **axes** –Axes that starts and ends apply to.
- **step** –Slice step, step = 1 if step is not specified

**Returns** Output tensor pointer, created by this slice function

template<typename T>

*TensorBase* \***pad** (T \*input\_element, const std::vector<int> &input\_shape, const std::vector<int> &pads, const padding\_mode\_t mode, *TensorBase* \*const\_value = nullptr)

Pad input tensor.

**Parameters**

- **input\_element** –Data pointer of input tensor
- **input\_shape** –Shape of input tensor
- **pads** –The number of padding elements to add, pads format should be: [x1\_begin, x2\_begin, ..., x1\_end, x2\_end, ...]
- **mode** –Supported modes: constant(default), reflect, edge
- **const\_value** –(Optional) A scalar value to be used if the mode chosen is constant

**Returns** Output tensor pointer

*TensorBase* \***pad** (*TensorBase* \*input, const std::vector<int> &pads, const padding\_mode\_t mode, *TensorBase* \*const\_value = nullptr)

Pad input tensor.

**Parameters**

- **input** –Input tensor pointer
- **pads** –Padding elements to add, pads format should be: [x1\_begin, x2\_begin, ..., x1\_end, x2\_end, ...]
- **mode** –Supported modes: constant(default), reflect, edge
- **const\_value** –(Optional) A scalar value to be used if the mode chosen is constant

**Returns** Output tensor pointer

template<typename T>

bool **compare\_elements** (const T \*gt\_elements, float epsilon = 1e-6, bool verbose = false)

Compare the elements of two Tensor.

**Parameters**

- **gt\_elements** –The ground truth elements
- **epsilon** –The max error of two element

- **verbose** –If true, print the detail of results

**Returns** true if all elements are equal otherwise false

int **get\_element\_index** (const std::vector<int> &axis\_index)

Get the index of element.

**Parameters** **axis\_index** –The coordinates of element

**Returns** int the index of element

std::vector<int> **get\_element\_coordinates** (int index)

Get the coordinates of element.

**Parameters** **index** –The index of element

**Returns** The coordinates of element

template<typename T>

T **get\_element** (int index)

Get a element of Tensor by index.

**Parameters** **index** –The index of element

**Returns** The element of tensor

template<typename T>

T **get\_element** (const std::vector<int> &axis\_index)

Get a element of Tensor.

**Parameters** **axis\_index** –The index of element

**Returns** The element of tensor

size\_t **set\_preload\_addr** (void \*addr, size\_t size)

Set preload address of Tensor.

**Parameters**

- **addr** –The address of preload data
- **size** –Size of preload data

**Returns** The size of preload data

inline virtual void **preload** ()

Preload the data of Tensor.

void **reset\_bias\_layout** (quant\_type\_t op\_quant\_type, bool is\_depthwise)

Reset the layout of Tensor.

**Warning:** Only available for Convolution. Don't use it unless you know exactly what it does.

**Parameters**

- **op\_quant\_type** –The quant type of operation
- **is\_depthwise** –Whether is depthwise convolution

virtual void **print** (bool print\_data = false)

print the information of *TensorBase*

**Parameters** **print\_data** –Whether print the data

## Public Members

int **size**

size of element including padding

std::vector<int> **shape**

shape of Tensor

dtype\_t **dtype**

data type of element

int **exponent**

exponent of element

bool **auto\_free**

free element when object destroy

std::vector<int> **axis\_offset**

element offset of each axis

void \***data**

data pointer

void \***cache**

cache pointer , used for preload and do not need to free

uint32\_t **caps**

flags indicating the type of memory

### Public Static Functions

```
static void slice (TensorBase *input, TensorBase *output, const std::vector<int> &start, const
                 std::vector<int> &end, const std::vector<int> &axes = {}, const std::vector<int> &step =
                 {})
```

Produces a slice along multiple axes.

|   |
|---|
| <b>Warning:</b> The length of start, end and step must be same as the shape of input tensor |
|---|

#### Parameters

- **input** –Input Tensor
- **output** –Output Tensor
- **start** –Starting indicesd
- **end** –Ending indices
- **axes** –Axes that starts and ends apply to.
- **step** –Slice step, step = 1 if step is not specified

**Returns** Void

## 4.2 Module API Reference

The Module is the base class for operators in esp-dl, and all operators inherit from this base class. This base class defines the basic interfaces for operators, enabling the model layer to automatically execute operators and manage memory planning.

## 4.2.1 Header File

- [esp-dl/dl/module/include/dl\\_module\\_base.hpp](#)

## 4.2.2 Classes

class **Module**

Base class for module.

### Public Functions

**Module** (const char \*name = NULL, module\_inplace\_t inplace = MODULE\_NON\_INPLACE, quant\_type\_t quant\_type = QUANT\_TYPE\_NONE)

Construct a new *Module* object.

#### Parameters

- **name** –Name of module.
- **inplace** –Inplace operation mode
- **quant\_type** –Quantization type

virtual **~Module** ()

Destroy the *Module* object. Return resource.

virtual std::vector<std::vector<int>> **retrieve\_inputs\_shape** (std::vector<std::vector<int>> &input\_shapes, std::vector<dl::TensorBase\*> inputs = {})

Retrieve the shape of this module' s inputs.

#### Parameters

- **input\_shapes** –The feature\_map shape of this module' s inputs.
- **inputs** –If the module has constant inputs, the order and quantity of the parameters passed must be consistent with those defined in ONNX. If it is a constant input, pass in its *TensorBase* pointer; if not, pass in nullptr.

**Returns** std::vector<std::vector<int>> Input shapes

virtual std::vector<TensorBase\*> **retrieve\_inputs** (std::vector<TensorBase\*> &tensors, std::vector<dl::TensorBase\*> inputs = {})

Retrieve the module' s inputs.

#### Parameters

- **tensors** –All inputs and outputs from MemoryManager
- **inputs** –If the module has constant inputs, the order and quantity of the parameters passed must be consistent with those defined in ONNX. If it is a constant input, pass in its *TensorBase* pointer; if not, pass in nullptr.

**Returns** std::vector<TensorBase \*> The final inputs.

inline virtual std::vector<int> **get\_outputs\_index** ()

Get the tensor index of this module' s outputs.

**Returns** Tensor index of model' s tensors

virtual std::vector<std::vector<int>> **get\_output\_shape** (std::vector<std::vector<int>> &input\_shapes) = 0

Calculate output shape by input shape.

**Parameters** **input\_shapes** –Input shapes

**Returns** outputs shapes



virtual void **forward** (std::vector<dl::TensorBase\*> &tensors, runtime\_mode\_t mode =  
 RUNTIME\_MODE\_AUTO) = 0

Run the module, high-level interface for model layer.

**Parameters**

- **tensors** –All inputs and outputs from MemoryManager
- **mode** –Runtime mode, default is RUNTIME\_MODE\_AUTO

inline virtual void **forward\_args** (void \*args)

Run the module, Low-level interface for base layer and multi-core processing.

**Parameters** **args** –ArgsType, arithArgsType, resizeArgsType and so on

inline virtual void **print** ()

print module information

inline virtual void **set\_preload\_addr** (void \*addr, size\_t size)

set preload RAM pointer

**Parameters**

- **addr** –Internal RAM address, should be aligned to 16 bytes
- **size** –The size of RAM address

inline virtual void **preload** ()

Perform a preload operation.

|                                 |
|---------------------------------|
| <b>Warning:</b> Not implemented |
|---------------------------------|

inline virtual void **reset** ()

reset all state of module, include inputs, outputs and preload cache setting

virtual void **run** (TensorBase \*input, TensorBase \*output, runtime\_mode\_t mode =  
 RUNTIME\_MODE\_AUTO)

Run the module with single input and single output.

**Parameters**

- **input** –Input tensor
- **output** –Output tensor
- **mode** –Runtime mode

virtual void **run** (std::vector<dl::TensorBase\*> inputs, std::vector<dl::TensorBase\*> outputs,  
 runtime\_mode\_t mode = RUNTIME\_MODE\_AUTO)

Run the module by inputs and outputs.

**Parameters**

- **inputs** –Input tensors
- **outputs** –Output tensors
- **mode** –Runtime mode

inline virtual void **get\_param\_memory\_size** (mem\_info \*in\_fbs, mem\_info \*out\_fbs, fbs::FbsModel  
 \*fbs\_model)

Get the memory size of module parameters.

**Parameters**

- **in\_fbs** –Memory info of module parameters inside Flatbuffers model
- **out\_fbs** –Memory info of module parameters outside Flatbuffers model
- **fbs\_model** –Flatbuffers model

## Public Members

char **\*name**

Name of module.

module\_inplace\_t **inplace**

Inplace type.

quant\_type\_t **quant\_type**

Quantization type.

std::vector<int> **m\_inputs\_index**

Tensor index of model' s tensors that used for inputs.

std::vector<int> **m\_outputs\_index**

Tensor index of model' s tensors that used for outputs.

## Public Static Functions

static inline *Module* \***deserialize** (fbs::*FbsModel* \*fbs\_model, std::string node\_name)  
create module instance by node serialization information

### Parameters

- **fbs\_model** –Flatbuffer' s model
- **node\_name** –The node name in model' s graph

**Returns** The pointer of module instance

static inline void **get\_param\_memory\_size** (dl::*TensorBase* \*param, mem\_info \*in\_fbs, mem\_info \*out\_fbs, fbs::*FbsModel* \*fbs\_model)

Get the memory size of module parameters.

### Parameters

- **param** –*Module* parameter tensor
- **in\_fbs** –Memory info of parameter inside Flatbuffers model
- **out\_fbs** –Memory info of parameter outside Flatbuffers model
- **fbs\_model** –Flatbuffers model

static inline void **get\_param\_memory\_size** (const std::vector<dl::*TensorBase*\*> &params, mem\_info \*in\_fbs, mem\_info \*out\_fbs, fbs::*FbsModel* \*fbs\_model)

Get the memory size of module parameters.

### Parameters

- **params** –*Module* parameter tensors
- **in\_fbs** –Memory info of parameters inside Flatbuffers model
- **out\_fbs** –Memory info of parameters outside Flatbuffers model
- **fbs\_model** –Flatbuffers model

## 4.2.3 Header File

- [esp-dl/dl/module/include/dl\\_module\\_creator.hpp](#)

## 4.2.4 Classes

class **ModuleCreator**

Singleton class for registering modules.

### Public Types

using **Creator** = std::function<*Module*\*(fbs::FbsModel\*, std::string)>

*Module* creator function type.

### Public Functions

inline void **register\_module** (const std::string &op\_type, *Creator* creator)

Register a module creator to the module creator map This function allows for the dynamic registration of new module types and their corresponding creator functions at runtime. By associating the module type name with the creator function, the system can flexibly create instances of various modules.

#### Parameters

- **op\_type** –The module type name, used as the key in the map
- **creator** –The module creator function, used to create modules of a specific type

inline *Module* \***create** (fbs::FbsModel \*fbs\_model, const std::string &op\_type, const std::string name)

Create module instance pointer.

#### Parameters

- **fbs\_model** –Flatbuffer model pointer
- **op\_type** –Module/Operator type
- **name** –*Module* name

**Returns** *Module* instance pointer

inline void **register\_dl\_modules** ()

Pre-register the already implemented modules.

inline void **print** ()

Print all modules has been registered.

inline void **clear** ()

Clear all modules has been registered.

### Public Static Functions

static inline *ModuleCreator* \***get\_instance** ()

Get instance of *ModuleCreator* by this function. It is only safe method to get instance of *ModuleCreator* because *ModuleCreator* is a singleton class.

**Returns** *ModuleCreator* instance pointer

## 4.3 Model API Reference

This section covers model loading and static memory planning, making it convenient for users to directly load and run ESPDL models.

### 4.3.1 Header File

- [esp-dl/dl/model/include/dl\\_model\\_base.hpp](#)

### 4.3.2 Macros

`DL_LOG_INFER_LATENCY_INIT_WITH_SIZE` (size)  
`DL_LOG_INFER_LATENCY_INIT` ()  
`DL_LOG_INFER_LATENCY_START` ()  
`DL_LOG_INFER_LATENCY_END` ()  
`DL_LOG_INFER_LATENCY_PRINT` (prefix, key)  
`DL_LOG_INFER_LATENCY_END_PRINT` (prefix, key)  
`DL_LOG_INFER_LATENCY_ARRAY_INIT_WITH_SIZE` (n, size)  
`DL_LOG_INFER_LATENCY_ARRAY_INIT` (n)  
`DL_LOG_INFER_LATENCY_ARRAY_START` (i)  
`DL_LOG_INFER_LATENCY_ARRAY_END` (i)  
`DL_LOG_INFER_LATENCY_ARRAY_PRINT` (i, prefix, key)  
`DL_LOG_INFER_LATENCY_ARRAY_END_PRINT` (i, prefix, key)

### 4.3.3 Classes

class **Model**

Neural Network *Model*.

#### Public Functions

**Model** (const char \*rodata\_address\_or\_partition\_label\_or\_path, fbs::model\_location\_type\_t location = fbs::MODEL\_LOCATION\_IN\_FLASH\_RODATA, int max\_internal\_size = 0, memory\_manager\_t mm\_type = MEMORY\_MANAGER\_GREEDY, uint8\_t \*key = nullptr, bool param\_copy = true)

Create the *Model* object by rodata address or partition label.

#### Parameters

- **rodata\_address\_or\_partition\_label\_or\_path** –The address of model data while location is `MODEL_LOCATION_IN_FLASH_RODATA`. The label of partition while location is `MODEL_LOCATION_IN_FLASH_PARTITION`. The path of model while location is `MODEL_LOCATION_IN_SDCARD`.
- **location** –The model location.
- **max\_internal\_size** –In bytes. Limit the max internal size usage. Only take effect when there's a PSRAM, and you want to alloc memory on internal RAM first.
- **mm\_type** –Type of memory manager
- **key** –The key of encrypted model.

- **param\_copy** –Set to false to avoid copy model parameters from FLASH to PSRAM. Only set this param to false when your PSRAM resource is very tight. This saves PSRAM and sacrifices the performance of model inference because the frequency of PSRAM is higher than FLASH. Only takes effect when MODEL\_LOCATION\_IN\_FLASH\_RODATA(CONFIG\_SPIRAM\_RODATA not set) or MODEL\_LOCATION\_IN\_FLASH\_PARTITION.

**Model** (const char \*rodata\_address\_or\_partition\_label\_or\_path, int model\_index, fbs::model\_location\_type\_t location = fbs::MODEL\_LOCATION\_IN\_FLASH\_RODATA, int max\_internal\_size = 0, memory\_manager\_t mm\_type = MEMORY\_MANAGER\_GREEDY, uint8\_t \*key = nullptr, bool param\_copy = true)

Create the *Model* object by rodata address or partition label.

#### Parameters

- **rodata\_address\_or\_partition\_label\_or\_path** –The address of model data while location is MODEL\_LOCATION\_IN\_FLASH\_RODATA. The label of partition while location is MODEL\_LOCATION\_IN\_FLASH\_PARTITION. The path of model while location is MODEL\_LOCATION\_IN\_SDCARD.
- **model\_index** –The model index of packed models.
- **location** –The model location.
- **max\_internal\_size** –In bytes. Limit the max internal size usage. Only take effect when there's a PSRAM, and you want to alloc memory on internal RAM first.
- **mm\_type** –Type of memory manager
- **key** –The key of encrypted model.
- **param\_copy** –Set to false to avoid copy model parameters from FLASH to PSRAM. Only set this param to false when your PSRAM resource is very tight. This saves PSRAM and sacrifices the performance of model inference because the frequency of PSRAM is higher than FLASH. Only takes effect when MODEL\_LOCATION\_IN\_FLASH\_RODATA(CONFIG\_SPIRAM\_RODATA not set) or MODEL\_LOCATION\_IN\_FLASH\_PARTITION.

**Model** (const char \*rodata\_address\_or\_partition\_label\_or\_path, const char \*model\_name, fbs::model\_location\_type\_t location = fbs::MODEL\_LOCATION\_IN\_FLASH\_RODATA, int max\_internal\_size = 0, memory\_manager\_t mm\_type = MEMORY\_MANAGER\_GREEDY, uint8\_t \*key = nullptr, bool param\_copy = true)

Create the *Model* object by rodata address or partition label.

#### Parameters

- **rodata\_address\_or\_partition\_label\_or\_path** –The address of model data while location is MODEL\_LOCATION\_IN\_FLASH\_RODATA. The label of partition while location is MODEL\_LOCATION\_IN\_FLASH\_PARTITION. The path of model while location is MODEL\_LOCATION\_IN\_SDCARD.
- **model\_name** –The model name of packed models.
- **location** –The model location.
- **max\_internal\_size** –In bytes. Limit the max internal size usage. Only take effect when there's a PSRAM, and you want to alloc memory on internal RAM first.
- **mm\_type** –Type of memory manager
- **key** –The key of encrypted model.
- **param\_copy** –Set to false to avoid copy model parameters from FLASH to PSRAM. Only set this param to false when your PSRAM resource is very tight. This saves PSRAM and sacrifices the performance of model inference because the frequency of PSRAM is higher than FLASH. Only takes effect when MODEL\_LOCATION\_IN\_FLASH\_RODATA(CONFIG\_SPIRAM\_RODATA not set) or MODEL\_LOCATION\_IN\_FLASH\_PARTITION.

**Model** (fbs::FbsModel \*fbs\_model, int internal\_size = 0, memory\_manager\_t mm\_type = MEMORY\_MANAGER\_GREEDY)

Create the *Model* object by fbs\_model.

#### Parameters

- **fbs\_model** –The fbs model.
- **internal\_size** –Internal ram size, in bytes
- **mm\_type** –Type of memory manager

virtual **~Model** ()

Destroy the *Model* object.

virtual esp\_err\_t **load** (const char \*rodata\_address\_or\_partition\_label\_or\_path, fbs::model\_location\_type\_t location = fbs::MODEL\_LOCATION\_IN\_FLASH\_RODATA, uint8\_t \*key = nullptr, bool param\_copy = true)

Load model graph and parameters from FLASH or sdcard.

#### Parameters

- **rodata\_address\_or\_partition\_label\_or\_path** –The address of model data while location is MODEL\_LOCATION\_IN\_FLASH\_RODATA. The label of partition while location is MODEL\_LOCATION\_IN\_FLASH\_PARTITION. The path of model while location is MODEL\_LOCATION\_IN\_SDCARD.
- **location** –The model location.
- **key** –The key of encrypted model.
- **param\_copy** –Set to false to avoid copy model parameters from FLASH to PSRAM. Only set this param to false when your PSRAM resource is very tight. This saves PSRAM and sacrifices the performance of model inference because the frequency of PSRAM is higher than FLASH. Only takes effect when MODEL\_LOCATION\_IN\_FLASH\_RODATA(CONFIG\_SPIRAM\_RODATA not set) or MODEL\_LOCATION\_IN\_FLASH\_PARTITION.

#### Returns

- ESP\_OK Success
- ESP\_FAIL Failed

virtual esp\_err\_t **load** (const char \*rodata\_address\_or\_partition\_label\_or\_path, fbs::model\_location\_type\_t location = fbs::MODEL\_LOCATION\_IN\_FLASH\_RODATA, int model\_index = 0, uint8\_t \*key = nullptr, bool param\_copy = true)

Load model graph and parameters from FLASH or sdcard.

#### Parameters

- **rodata\_address\_or\_partition\_label\_or\_path** –The address of model data while location is MODEL\_LOCATION\_IN\_FLASH\_RODATA. The label of partition while location is MODEL\_LOCATION\_IN\_FLASH\_PARTITION. The path of model while location is MODEL\_LOCATION\_IN\_SDCARD.
- **location** –The model location.
- **model\_index** –The model index of packed models.
- **key** –The key of encrypted model.
- **param\_copy** –Set to false to avoid copy model parameters from FLASH to PSRAM. Only set this param to false when your PSRAM resource is very tight. This saves PSRAM and sacrifices the performance of model inference because the frequency of PSRAM is higher than FLASH. Only takes effect when MODEL\_LOCATION\_IN\_FLASH\_RODATA(CONFIG\_SPIRAM\_RODATA not set) or MODEL\_LOCATION\_IN\_FLASH\_PARTITION.

#### Returns

- ESP\_OK Success
- ESP\_FAIL Failed

virtual esp\_err\_t **load** (const char \*rodata\_address\_or\_partition\_label\_or\_path, fbs::model\_location\_type\_t location = fbs::MODEL\_LOCATION\_IN\_FLASH\_RODATA, const char \*model\_name = nullptr, uint8\_t \*key = nullptr, bool param\_copy = true)

Load model graph and parameters from FLASH or sdcard.

#### Parameters

- **rodata\_address\_or\_partition\_label\_or\_path** –The address of model data while location is MODEL\_LOCATION\_IN\_FLASH\_RODATA. The label of par-

tion while location is `MODEL_LOCATION_IN_FLASH_PARTITION`. The path of model while location is `MODEL_LOCATION_IN_SDCARD`.

- **location** –The model location.
- **model\_name** –The model name of packed models.
- **key** –The key of encrypted model.
- **param\_copy** –Set to false to avoid copy model parameters from FLASH to PSRAM. Only set this param to false when your PSRAM resource is very tight. This saves PSRAM and sacrifices the performance of model inference because the frequency of PSRAM is higher than FLASH. Only takes effect when `MODEL_LOCATION_IN_FLASH_RODATA(CONFIG_SPIRAM_RODATA not set)` or `MODEL_LOCATION_IN_FLASH_PARTITION`.

#### Returns

- `ESP_OK` Success
- `ESP_FAIL` Failed

virtual `esp_err_t load` (`fbs::FbsModel *fbs_model`)

Load model graph and parameters from Flatbuffers model.

**Parameters** `fbs_model` –The FlatBuffers model

#### Returns

- `ESP_OK` Success
- `ESP_FAIL` Failed

virtual void **build** (`size_t max_internal_size`, `memory_manager_t mm_type = MEMORY_MANAGER_GREEDY`, `bool preload = false`)

Allocate memory for the model.

#### Parameters

- **max\_internal\_size** –In bytes. Limit the max internal size usage. Only take effect when there' s a PSRAM, and you want to alloc memory on internal RAM first.
- **mm\_type** –Type of memory manager
- **preload** –Whether to preload the model' s parameters to internal ram (not implemented yet)

virtual void **run** (`runtime_mode_t mode = RUNTIME_MODE_SINGLE_CORE`)

Run the model module by module.

**Parameters** `mode` –Runtime mode.

virtual void **run** (`TensorBase *input`, `runtime_mode_t mode = RUNTIME_MODE_SINGLE_CORE`)

Run the model module by module.

#### Parameters

- **input** –The model input.
- **mode** –Runtime mode.

virtual void **run** (`std::map<std::string, TensorBase*> &user_inputs`, `runtime_mode_t mode = RUNTIME_MODE_SINGLE_CORE`, `std::map<std::string, TensorBase*> user_outputs = {}`)

Run the model module by module.

#### Parameters

- **user\_inputs** –The model inputs.
- **mode** –Runtime mode.
- **user\_outputs** –It' s for debug to pecify the output of the intermediate layer; Under normal use, there is no need to pass a value to this parameter. If no parameter is passed, the default is the graphical output, which can be obtained through `Model::get_outputs()`.

`esp_err_t test` ()

Test whether the model inference result is correct. The model should contain `test_inputs` and `test_outputs`. Enable `export_test_values` option in `esp-ppq` to use this api.

**Returns** `esp_err_t`

std::map<std::string, mem\_info> **get\_memory\_info** ()

Get memory info.

**Returns** Memory usage statistics on internal and PSRAM.

std::map<std::string, module\_info> **get\_module\_info** ()

Get module info.

**Returns** return Type and latency of each module.

void **print\_module\_info** (const std::map<std::string, module\_info> &info, bool  
sort\_module\_by\_latency = false)

Print the module info obtained by get\_module\_info function.

**Parameters**

- **info** –
- **sort\_module\_by\_latency** –

void **profile\_memory** ()

Print model memory summary.

void **profile\_module** (bool sort\_module\_by\_latency = false)

Print module info summary. (Name, Type, Latency)

**Parameters** **sort\_module\_by\_latency** –True The module is printed in latency decreasing sort. False The module is printed in ONNX topological sort.

void **profile** (bool sort\_module\_by\_latency = false)

Combination of profile\_memory & profile\_module.

**Parameters** **sort\_module\_by\_latency** –True The module is printed in latency decreasing sort. False The module is printed in ONNX topological sort.

virtual std::map<std::string, *TensorBase*> &**get\_inputs** ()

Get inputs of model.

**Returns** The map of model input' s name and *TensorBase*\*

virtual *TensorBase* \***get\_intermediate** (const std::string &name)

Get intermediate *TensorBase* of model.

---

**Note:** When using memory manager, the content of *TensorBase*' s data may be overwritten by the outputs of other

---

**Parameters** **name** –The name of intermediate Tensor. operators.

**Returns** The intermediate *TensorBase*\*

virtual std::map<std::string, *TensorBase*> &**get\_outputs** ()

Get outputs of model.

**Returns** The map of model output' s name and *TensorBase*\*

virtual void **print** ()

Print the model.

inline virtual *FbsModel* \***get\_fbs\_model** ()

Get the fbs model instance.

**Returns** *fbs::FbsModel* \*



### 4.3.4 Header File

- `esp-dl/dl/model/include/dl_memory_manager.hpp`

### 4.3.5 Classes

#### class **MemoryManagerBase**

Memory manager base class, each model has its own memory manager TODO: share memory manager with different models.

Subclassed by `dl::memory::MemoryManagerGreedy`

#### Public Functions

inline **MemoryManagerBase** (int alignment = 16)

Construct a new Memory Manager Base object.

**Parameters** `alignment` –Memory address alignment

inline virtual **~MemoryManagerBase** ()

Destroy the MemoryManager object. Return resource.

inline virtual void **alloc** (fbs::FbsModel \*fbs\_model, std::vector<dl::module::Module\*> &execution\_plan)

Allocate memory for each tensor, include all input and output tensors.

#### Parameters

- `fbs_model` –FlatBuffer' s *Model*
- `execution_plan` –Topological sorted module list

inline virtual void **set\_preload\_addr** (std::vector<dl::module::Module\*> execution\_plan)

Set preload address for module' s parameters.

**Parameters** `execution_plan` –Topological sorted module list

virtual void **reset** ()

Reset the memory manager, free all memory for each tensor, include all input and output tensors.

*TensorBase* \***get\_tensor** (int index)

Get tensor by index.

**Parameters** `index` –The tensor index, type: int

**Returns** The *TensorBase* pointer

*TensorBase* \***get\_tensor** (const std::string &name)

Get tensor by name.

**Parameters** `name` –The tensor name, type: std::string

**Returns** The *TensorBase* pointer

int **get\_tensor\_index** (const std::string &name)

Get tensor index by name.

**Parameters** `name` –The tensor name, type: std::string

**Returns** The *TensorBase* index

void **root\_free** ()

Free psram root and internal root pointer.

```
inline void *get_psram_root ()
```

Get PSRAM root pointer.

**Returns** void\*

```
inline void *get_internal_root ()
```

Get internal ram root pointer.

**Returns** void\*

```
inline size_t get_psram_size ()
```

Get PSRAM size allocated by memory manager.

**Returns** size\_t

```
inline size_t get_internal_size ()
```

Get internal RAM size allocated by memory manager.

**Returns** size\_t

### Public Members

```
std::vector<TensorBase*> tensors
```

All tensors in the model

```
int alignment
```

The root pointer needs to be aligned must be a power of two

```
std::map<std::string, int> name2index
```

Tensor name to index map

```
size_t internal_size
```

The bytes of internal ram

```
size_t psram_size
```

The bytes of psram

```
class TensorInfo
```

Tensor info, include tensor name, shape, dtype, size, time range and call times, which is used to plan model memory.

### Public Functions

```
TensorInfo (std::string &name, int time_begin, int time_end, std::vector<int> shape, dtype_t dtype, int exponent, bool is_internal = false)
```

Construct a new Tensor Info object.

#### Parameters

- **name** –Tensor name
- **time\_begin** –Tensor lifetime begin
- **time\_end** –Tensor lifetime end
- **shape** –Tensor shape
- **dtype** –Tensor dtype
- **exponent** –Tensor exponent
- **is\_internal** –Is tensor in internal RAM or not

inline **~TensorInfo** ()

Destroy the Tensor Info object.

void **set\_inplace\_leader\_tensor** (*TensorInfo* \*tensor)

Set the inplace leader tensor object.

**Parameters** **tensor** –Inplace leader tensor

inline void **set\_inplace\_follower\_tensor** (*TensorInfo* \*tensor)

Set the inplace follower tensor object.

**Parameters** **tensor** –Inplace follower tensor

inline *TensorInfo* \***get\_inplace\_follower\_tensor** ()

Get the inplace follower tensor object.

**Returns** *TensorInfo*\* Inplace follower tensor

void **update\_time** (int new\_time)

Update Tensor lifetime.

**Parameters** **new\_time** –new tensor lifetime

*TensorBase* \***create\_tensor** (void \*internal\_root, void \*psram\_root)

Create a *TensorBase* object according to *TensorInfo*.

**Parameters**

- **internal\_root** –Internal RAM root pointer
- **psram\_root** –PSRAM root pointer

**Returns** *TensorBase*\*

inline bool **is\_inplaced** ()

Is inplaced or not.

**Returns** true if inplaced else false

inline uint32\_t **get\_offset** ()

Get the tensor offset.

**Returns** uint32\_t

inline void **set\_offset** (uint32\_t offset)

Set the tensor offset.

**Parameters** **offset** –

inline uint32\_t **get\_internal\_offset** ()

Get the internal offset.

**Returns** uint32\_t

inline bool **get\_internal\_state** ()

Get the internal state.

**Returns** true if is internal else false

inline void **set\_internal\_state** (bool is\_internal)

Set the internal state.

**Parameters** **is\_internal** –

inline void **set\_internal\_offset** (uint32\_t offset)

Set the internal offset.

**Parameters** **offset** –

```
inline int get_time_end ()  
    Get the lifetime end.  
    Returns int  
inline int get_time_begin ()  
    Get the lifetime begin.  
    Returns int  
inline size_t get_size ()  
    Get the tensor size.  
    Returns size_t  
inline std::string get_name ()  
    Get the tensor name.  
    Returns std::string  
inline std::vector<int> get_shape ()  
    Get the tensor shape.  
    Returns std::vector<int>  
inline void print ()  
    print tensor info
```

class **MemoryChunk**

Memory chunk, include size, is free, offset, alignment and tensor, which is used to simulate memory allocation.

### Public Functions

**MemoryChunk** (size\_t size, int is\_free, int alignment = 16)

Construct a new Memory Chunk object.

#### Parameters

- **size** –Memory chunk size
- **is\_free** –Whether free or not
- **alignment** –Memory chunk alignment

**MemoryChunk** (*TensorInfo* \*tensor, int alignment = 16)

Construct a new Memory Chunk object.

#### Parameters

- **tensor** –*TensorInfo*
- **alignment** –Memory chunk alignment

inline **~MemoryChunk** ()

Destroy the Memory Chunk object.

*MemoryChunk* \***merge\_free\_chunk** (*MemoryChunk* \*chunk)

Merge continuous free chunk.

**Parameters** **chunk** –

**Returns** *MemoryChunk*\*

*MemoryChunk* \***insert** (*TensorInfo* \*tensor)

Insert tensor into free chunk.

**Parameters** **tensor** –

**Returns** *MemoryChunk*\*

*MemoryChunk* \***extend** (*TensorInfo* \*tensor)

Extend free chunk and insert tensor.

**Parameters** *tensor* –

**Returns** *MemoryChunk*\*

inline void **free** ()

Free memory chunk, set *is\_free* to true and set *tensor* to nullptr.

size\_t **get\_aligned\_size** (size\_t size)

get aligned size, which is 16/alignemt bytes aligned

**Parameters** *size* –

**Returns** size\_t

### Public Members

size\_t **size**

Memeory chunk size

bool **is\_free**

Whether memory chunk is free or not

int **offset**

Offset relative to root pointer

int **alignment**

Memory address alignment

*TensorInfo* \***tensor**

Info of the tensor which occupies the memory

### 4.3.6 Header File

- [esp-dl/dl/model/include/dl\\_memory\\_manager\\_greedy.hpp](#)

### 4.3.7 Classes

class **MemoryManagerGreedy** : public dl::memory::MemoryManagerBase

Greedy memory manager, allocate memory for each tensor in the order of the execution plan.

#### Public Functions

inline **MemoryManagerGreedy** (int max\_internal\_size, int alignment = 16)

Construct a new Memory Manager Greedy object.

**Parameters**

- **max\_internal\_size** –In bytes. Limit the max internal size usage. Only take effect when there' s a PSRAM, and you want to alloc memory on internal RAM first
- **alignment** –Memory address alignment

inline **~MemoryManagerGreedy** ()

Destroy the Memory Manager Greedy object.

virtual void **alloc** (fbs::FbsModel \*fbs\_model, std::vector<dl::module::Module\*> &execution\_plan)

Allocate memory for each tensor, include all input and output tensors.

**Parameters**

- **fbs\_model** –FlatBuffer’ s *Model*
- **execution\_plan** –Topological sorted module list

virtual void **set\_preload\_addr** (std::vector<dl::module::Module\*> execution\_plan)

Set preload address for module’ s parameters.

**Parameters** **execution\_plan** –Topological sorted module list

void **free** ()

Free memory, include all tensors and memory list.

## 4.4 Fbs API Reference

The esp-dl model utilizes FlatBuffers to store information about parameters and the computation graph. Taking into account the encryption requirements of some models, this part has not been open-sourced. However, we provide a set of APIs to facilitate users in loading and parsing esp-dl models.

### 4.4.1 Header File

- [esp-dl/fbs\\_loader/include/fbs\\_loader.hpp](#)

### 4.4.2 Classes

class **FbsLoader**

Class for parser the flatbuffers.

#### Public Functions

**FbsLoader** (const char \*rodata\_address\_or\_partition\_label\_or\_path = nullptr, model\_location\_type\_t location = MODEL\_LOCATION\_IN\_FLASH\_RODATA)

Construct a new *FbsLoader* object.

**Parameters**

- **rodata\_address\_or\_partition\_label\_or\_path** –The address of model data while location is MODEL\_LOCATION\_IN\_FLASH\_RODATA. The label of partition while location is MODEL\_LOCATION\_IN\_FLASH\_PARTITION. The path of model while location is MODEL\_LOCATION\_IN\_SDCARD.
- **location** –The model location.

**~FbsLoader** ()

Destroy the *FbsLoader* object.

*FbsModel* \***load** (const uint8\_t \*key = nullptr, bool param\_copy = true)

Load the model. If there are multiple sub-models, the first sub-model will be loaded.

**Parameters**

- **key** –NULL or a 128-bit AES key, like {0x00, 0x01, 0x02, 0x03, 0x04, 0x05, 0x06, 0x07, 0x08, 0x09, 0x0a, 0x0b, 0x0c, 0x0d, 0x0e, 0x0f }

- **param\_copy** –Set to false to avoid copy model parameters from FLASH to PSRAM. Only set this param to false when your PSRAM resource is very tight. This saves PSRAM and sacrifices the performance of model inference because the frequency of PSRAM is higher than FLASH. Only takes effect when MODEL\_LOCATION\_IN\_FLASH\_RODATA(CONFIG\_SPIRAM\_RODATA not set) or MODEL\_LOCATION\_IN\_FLASH\_PARTITION.

**Returns** Return nullptr if loading fails. Otherwise return the pointer of *FbsModel*.

*FbsModel* \***load** (const int model\_index, const uint8\_t \*key = nullptr, bool param\_copy = true)

Load the model by model index.

#### Parameters

- **model\_index** –The index of model.
- **key** –NULL or a 128-bit AES key, like {0x00, 0x01, 0x02, 0x03, 0x04, 0x05, 0x06, 0x07, 0x08, 0x09, 0x0a, 0x0b, 0x0c, 0x0d, 0x0e, 0x0f}.
- **param\_copy** –Set to false to avoid copy model parameters from FLASH to PSRAM. Only set this param to false when your PSRAM resource is very tight. This saves PSRAM and sacrifices the performance of model inference because the frequency of PSRAM is higher than FLASH. Only takes effect when MODEL\_LOCATION\_IN\_FLASH\_RODATA(CONFIG\_SPIRAM\_RODATA not set) or MODEL\_LOCATION\_IN\_FLASH\_PARTITION.

**Returns** Return nullptr if loading fails. Otherwise return the pointer of *FbsModel*.

*FbsModel* \***load** (const char \*model\_name, const uint8\_t \*key = nullptr, bool param\_copy = true)

Load the model by model name.

#### Parameters

- **model\_name** –The name of model.
- **key** –NULL or a 128-bit AES key, like {0x00, 0x01, 0x02, 0x03, 0x04, 0x05, 0x06, 0x07, 0x08, 0x09, 0x0a, 0x0b, 0x0c, 0x0d, 0x0e, 0x0f}
- **param\_copy** –Set to false to avoid copy model parameters from FLASH to PSRAM. Only set this param to false when your PSRAM resource is very tight. This saves PSRAM and sacrifices the performance of model inference because the frequency of PSRAM is higher than FLASH. Only takes effect when MODEL\_LOCATION\_IN\_FLASH\_RODATA(CONFIG\_SPIRAM\_RODATA not set) or MODEL\_LOCATION\_IN\_FLASH\_PARTITION.

**Returns** Return nullptr if loading fails. Otherwise return the pointer of *FbsModel*.

int **get\_model\_num** ()

Get the number of models.

**Returns** The number of models

void **list\_models** ()

List all model' s name.

### 4.4.3 Header File

- esp-dl/fbs\_loader/include/fbs\_model.hpp

### 4.4.4 Classes

class **FbsModel**

Flatbuffer model object.

#### Public Functions

**FbsModel** (const void \*data, size\_t size, model\_location\_type\_t location, bool encrypt, bool rodata\_move, bool auto\_free, bool param\_copy)

Construct a new *FbsModel* object.

**Parameters**

- **data** –The data of model flatbuffers.
- **size** –The size of model flatbuffers in bytes.
- **location** –The location of model flatbuffers.
- **encrypt** –Whether the model flatbuffers is encrypted or not.
- **rodata\_move** –Whether the model flatbuffers is moved from FLASH rodata to PSRAM.
- **auto\_free** –Whether to free the model flatbuffers data when destroy this class instance.
- **param\_copy** –Whether to copy the parameter in flatbuffers.

**~FbsModel** ()

Destroy the *FbsModel* object.

void **print** ()

Print the model information.

std::vector<std::string> **topological\_sort** ()

Return vector of node name in the order of execution.

**Returns** topological sort of node name.

esp\_err\_t **get\_operation\_attribute** (std::string node\_name, std::string attribute\_name, int &ret\_value)

Get the attribute of node.

**Parameters**

- **node\_name** –The name of operation.
- **attribute\_name** –The name of attribute.
- **ret\_value** –The attribute value.

**Returns** esp\_err\_t Return ESP\_OK if get successfully. Otherwise return ESP\_FAIL.

esp\_err\_t **get\_operation\_attribute** (std::string node\_name, std::string attribute\_name, float &ret\_value)

Get the attribute of node.

**Parameters**

- **node\_name** –The name of operation.
- **attribute\_name** –The name of attribute.
- **ret\_value** –The attribute value.

**Returns** esp\_err\_t Return ESP\_OK if get successfully. Otherwise return ESP\_FAIL.

esp\_err\_t **get\_operation\_attribute** (std::string node\_name, std::string attribute\_name, std::string &ret\_value)

Get the attribute of node.

**Parameters**

- **node\_name** –The name of operation.
- **attribute\_name** –The name of attribute.
- **ret\_value** –The attribute value.

**Returns** esp\_err\_t Return ESP\_OK if get successfully. Otherwise return ESP\_FAIL.

esp\_err\_t **get\_operation\_attribute** (std::string node\_name, std::string attribute\_name, std::vector<int> &ret\_value)

Get the attribute of node.

**Parameters**

- **node\_name** –The name of operation.



- **attribute\_name** –The name of attribute.
- **ret\_value** –The attribute value.

**Returns** esp\_err\_t Return ESP\_OK if get successfully. Otherwise return ESP\_FAIL.

esp\_err\_t **get\_operation\_attribute** (std::string node\_name, std::string attribute\_name, std::vector<float> &ret\_value)

Get the attribute of node.

#### Parameters

- **node\_name** –The name of operation.
- **attribute\_name** –The name of attribute.
- **ret\_value** –The attribute value.

**Returns** esp\_err\_t Return ESP\_OK if get successfully. Otherwise return ESP\_FAIL.

esp\_err\_t **get\_operation\_attribute** (std::string node\_name, std::string attribute\_name, dl::quant\_type\_t &ret\_value)

Get the attribute of node.

#### Parameters

- **node\_name** –The name of operation.
- **attribute\_name** –The name of attribute.
- **ret\_value** –The attribute value.

**Returns** esp\_err\_t Return ESP\_OK if get successfully. Otherwise return ESP\_FAIL.

esp\_err\_t **get\_operation\_attribute** (std::string node\_name, std::string attribute\_name, dl::activation\_type\_t &ret\_value)

Get the attribute of node.

#### Parameters

- **node\_name** –The name of operation.
- **attribute\_name** –The name of attribute.
- **ret\_value** –The attribute value.

**Returns** esp\_err\_t Return ESP\_OK if get successfully. Otherwise return ESP\_FAIL.

esp\_err\_t **get\_operation\_attribute** (std::string node\_name, std::string attribute\_name, dl::resize\_mode\_t &ret\_value)

Get the attribute of node.

#### Parameters

- **node\_name** –The name of operation.
- **attribute\_name** –The name of attribute.
- **ret\_value** –The attribute value.

**Returns** esp\_err\_t Return ESP\_OK if get successfully. Otherwise return ESP\_FAIL.

esp\_err\_t **get\_operation\_attribute** (std::string node\_name, std::string attribute\_name, dl::TensorBase \* &ret\_value)

Get the attribute of node.

#### Parameters

- **node\_name** –The name of operation.
- **attribute\_name** –The name of attribute.
- **ret\_value** –The attribute value.

**Returns** esp\_err\_t Return ESP\_OK if get successfully. Otherwise return ESP\_FAIL.

esp\_err\_t **get\_operation\_output\_shape** (std::string node\_name, int index, std::vector<int> &ret\_value)

Get operation output shape.

#### Parameters

- **node\_name** –The name of operation.
- **index** –The index of outputs
- **ret\_value** –Return shape value.

**Returns** `esp_err_t` Return `ESP_OK` if get successfully. Otherwise return `ESP_FAIL`.

`esp_err_t` **get\_operation\_inputs\_and\_outputs** (`std::string` node\_name, `std::vector<std::string>` &inputs, `std::vector<std::string>` &outputs)

Get the attribute of node.

**Parameters**

- **node\_name** –The name of operation.
- **inputs** –The vector of operation inputs.
- **outputs** –The vector of operation outputs.

**Returns** `esp_err_t` Return `ESP_OK` if get successfully. Otherwise return `ESP_FAIL`.

`std::string` **get\_operation\_type** (`std::string` node\_name)

Get operation type, “Conv”, “Linear” etc.

**Parameters** **node\_name** –The name of operation

**Returns** The type of operation.

`dl::TensorBase*` **get\_operation\_parameter** (`std::string` node\_name, `int` index = 1, `uint32_t` caps = `MALLOC_CAP_DEFAULT`)

Return if the variable is a parameter.

**Parameters**

- **node\_name** –The name of operation
- **index** –The index of the variable
- **caps** –Bitwise OR of `MALLOC_CAP_*` flags indicating the type of memory to be returned

**Returns** `dl::TensorBase*`

`dl::TensorBase*` **get\_operation\_lut** (`std::string` node\_name, `uint32_t` caps = `MALLOC_CAP_DEFAULT`, `std::string` attribute\_name = “lut”)

Get LUT(Look Up Table) if the operation has LUT.

**Parameters**

- **node\_name** –The name of operation
- **caps** –Bitwise OR of `MALLOC_CAP_*` flags indicating the type of memory to be returned
- **attribute\_name** –The name of LUT attribute

**Returns** `dl::TensorBase*`

`bool` **is\_parameter** (`std::string` name)

return true if the variable is a parameter

**Parameters** **name** –Variable name

**Returns** true if the variable is a parameter else false

`const void*` **get\_tensor\_raw\_data** (`std::string` tensor\_name)

Get the raw data of `FlatBuffers::DL::Tensor`.

**Parameters** **tensor\_name** –The name of Tensor.

**Returns** `uint8_t*` The pointer of raw data.

`dl::dtype_t` **get\_tensor\_dtype** (`std::string` tensor\_name)

Get the element type of tensor tensor.

**Parameters** **tensor\_name** –The tensor name.

**Returns** `FlatBuffers::DL::TensorDataType`

`std::vector<int>` **get\_tensor\_shape** (`std::string` tensor\_name)

Get the shape of tensor.

**Parameters** **tensor\_name** –The name of tensor.

**Returns** `std::vector<int>` The shape of tensor.

std::vector<int> **get\_tensor\_exponents** (std::string tensor\_name)

Get the exponents of tensor.

**Warning:** When quantization is PER\_CHANNEL, the size of exponents is same as out\_channels. When quantization is PER\_TENSOR, the size of exponents is 1.

**Parameters** **tensor\_name** –The name of tensor.

**Returns** The exponents of tensor.

dl::dtype\_t **get\_value\_info\_dtype** (std::string var\_name)

Get the element type of value\_info.

**Parameters** **var\_name** –The value\_info name.

**Returns** dl::dtype\_t

std::vector<int> **get\_value\_info\_shape** (std::string var\_name)

Get the shape of value\_info.

**Parameters** **var\_name** –The value\_info name.

**Returns** the shape of value\_info.

int **get\_value\_info\_exponent** (std::string var\_name)

Get the exponent of value\_info. Only support PER\_TENSOR quantization.

**Parameters** **var\_name** –The value\_info name.

**Returns** the exponent of value\_info

const void \***get\_test\_input\_tensor\_raw\_data** (std::string tensor\_name)

Get the raw data of test input tensor.

**Parameters** **tensor\_name** –The name of test input tensor.

**Returns** uint8\_t \* The pointer of raw data.

const void \***get\_test\_output\_tensor\_raw\_data** (std::string tensor\_name)

Get the raw data of test output tensor.

**Parameters** **tensor\_name** –The name of test output tensor.

**Returns** uint8\_t \* The pointer of raw data.

dl::TensorBase \***get\_test\_input\_tensor** (std::string tensor\_name)

Get the test input tensor.

**Parameters** **tensor\_name** –The name of test input tensor.

**Returns** The pointer of tensor.

dl::TensorBase \***get\_test\_output\_tensor** (std::string tensor\_name)

Get the test output tensor.

**Parameters** **tensor\_name** –The name of test output tensor.

**Returns** The pointer of tensor.

std::vector<std::string> **get\_test\_outputs\_name** ()

Get the name of test outputs.

**Returns** the name of test outputs

std::vector<std::string> **get\_graph\_inputs** ()

Get the graph inputs.

**Returns** the name of inputs

std::vector<std::string> **get\_graph\_outputs** ()

Get the graph outputs.

**Returns** the name of outputs

void **clear\_map** ()

Clear all map.

void **load\_map** ()

Load all map.

std::string **get\_model\_name** ()

Get the model name.

**Returns** the name of model

int64\_t **get\_model\_version** ()

Get the model version.

**Returns** The version of model

std::string **get\_model\_doc\_string** ()

Get the model doc string.

**Returns** The doc string of model

void **get\_model\_size** (size\_t \*internal\_size, size\_t \*psram\_size, size\_t \*psram\_rodata\_size, size\_t \*flash\_size)

Get the model size.

**Parameters**

- **internal\_size** –Flatbuffers model internal RAM usage
- **psram\_size** –Flatbuffers model PSRAM usage
- **psram\_rodata\_size** –Flatbuffers model PSRAM rodata usage. If CONFIG\_SPIRAM\_RODATA option is on, \ Flatbuffers model in FLASH rodata will be copied to PSRAM
- **flash\_size** –Flatbuffers model FLASH usage

bool **memory\_addr\_in\_model** (void \*addr)

Whether the memory address within the fbs model or not.

**Parameters** **addr** –memory address

**Returns** true if memory address within the fbs model else false

## Public Members

bool **m\_param\_copy**

copy flatbuffers param or not.



# Index

## D

- dl::memory::MemoryChunk (C++ class), 57
- dl::memory::MemoryChunk::~~MemoryChunk (C++ function), 57
- dl::memory::MemoryChunk::alignment (C++ member), 58
- dl::memory::MemoryChunk::extend (C++ function), 57
- dl::memory::MemoryChunk::free (C++ function), 58
- dl::memory::MemoryChunk::get\_aligned\_size (C++ function), 58
- dl::memory::MemoryChunk::insert (C++ function), 57
- dl::memory::MemoryChunk::is\_free (C++ member), 58
- dl::memory::MemoryChunk::MemoryChunk (C++ function), 57
- dl::memory::MemoryChunk::merge\_free\_chunk (C++ function), 57
- dl::memory::MemoryChunk::offset (C++ member), 58
- dl::memory::MemoryChunk::size (C++ member), 58
- dl::memory::MemoryChunk::tensor (C++ member), 58
- dl::memory::MemoryManagerBase (C++ class), 54
- dl::memory::MemoryManagerBase::~~MemoryManagerBase (C++ function), 54
- dl::memory::MemoryManagerBase::alignment (C++ member), 55
- dl::memory::MemoryManagerBase::alloc (C++ function), 54
- dl::memory::MemoryManagerBase::get\_internal\_root (C++ function), 55
- dl::memory::MemoryManagerBase::get\_internal\_size (C++ function), 55
- dl::memory::MemoryManagerBase::get\_psrain\_root (C++ function), 54
- dl::memory::MemoryManagerBase::get\_psrain\_size (C++ function), 55
- dl::memory::MemoryManagerBase::get\_tensor (C++ function), 54
- dl::memory::MemoryManagerBase::get\_tensor\_index (C++ function), 54
- dl::memory::MemoryManagerBase::internal\_size (C++ member), 55
- dl::memory::MemoryManagerBase::MemoryManagerBase (C++ function), 54
- dl::memory::MemoryManagerBase::name2index (C++ member), 55
- dl::memory::MemoryManagerBase::psram\_size (C++ member), 55
- dl::memory::MemoryManagerBase::reset (C++ function), 54
- dl::memory::MemoryManagerBase::root\_free (C++ function), 54
- dl::memory::MemoryManagerBase::set\_preload\_addr (C++ function), 54
- dl::memory::MemoryManagerBase::tensors (C++ member), 55
- dl::memory::MemoryManagerGreedy (C++ class), 58
- dl::memory::MemoryManagerGreedy::~~MemoryManagerGreedy (C++ function), 58
- dl::memory::MemoryManagerGreedy::alloc (C++ function), 59
- dl::memory::MemoryManagerGreedy::free (C++ function), 59
- dl::memory::MemoryManagerGreedy::MemoryManagerGreedy (C++ function), 58
- dl::memory::MemoryManagerGreedy::set\_preload\_addr (C++ function), 59
- dl::memory::TensorInfo (C++ class), 55
- dl::memory::TensorInfo::~~TensorInfo (C++ function), 55
- dl::memory::TensorInfo::create\_tensor (C++ function), 56
- dl::memory::TensorInfo::get\_inplace\_follower\_tensors (C++ function), 56
- dl::memory::TensorInfo::get\_internal\_offset (C++ function), 56
- dl::memory::TensorInfo::get\_internal\_state (C++ function), 56
- dl::memory::TensorInfo::get\_name (C++ function), 57
- dl::memory::TensorInfo::get\_offset (C++ function), 56
- dl::memory::TensorInfo::get\_shape (C++ function), 57
- dl::memory::TensorInfo::get\_size (C++ function), 57
- dl::memory::TensorInfo::get\_time\_begin (C++ function), 57

- [\(C++ function\), 57](#)
- [dl::memory::TensorInfo::get\\_time\\_end \(C++ function\), 56](#)
- [dl::memory::TensorInfo::is\\_inplaced \(C++ function\), 56](#)
- [dl::memory::TensorInfo::print \(C++ function\), 57](#)
- [dl::memory::TensorInfo::set\\_inplace\\_follower \(C++ function\), 56](#)
- [dl::memory::TensorInfo::set\\_inplace\\_leader \(C++ function\), 56](#)
- [dl::memory::TensorInfo::set\\_internal\\_offset \(C++ function\), 56](#)
- [dl::memory::TensorInfo::set\\_internal\\_state \(C++ function\), 56](#)
- [dl::memory::TensorInfo::set\\_offset \(C++ function\), 56](#)
- [dl::memory::TensorInfo::TensorInfo \(C++ function\), 55](#)
- [dl::memory::TensorInfo::update\\_time \(C++ function\), 56](#)
- [dl::Model \(C++ class\), 49](#)
- [dl::Model::~~Model \(C++ function\), 51](#)
- [dl::Model::build \(C++ function\), 52](#)
- [dl::Model::get\\_fbs\\_model \(C++ function\), 53](#)
- [dl::Model::get\\_inputs \(C++ function\), 53](#)
- [dl::Model::get\\_intermediate \(C++ function\), 53](#)
- [dl::Model::get\\_memory\\_info \(C++ function\), 53](#)
- [dl::Model::get\\_module\\_info \(C++ function\), 53](#)
- [dl::Model::get\\_outputs \(C++ function\), 53](#)
- [dl::Model::load \(C++ function\), 51, 52](#)
- [dl::Model::Model \(C++ function\), 49, 50](#)
- [dl::Model::print \(C++ function\), 53](#)
- [dl::Model::print\\_module\\_info \(C++ function\), 53](#)
- [dl::Model::profile \(C++ function\), 53](#)
- [dl::Model::profile\\_memory \(C++ function\), 53](#)
- [dl::Model::profile\\_module \(C++ function\), 53](#)
- [dl::Model::run \(C++ function\), 52](#)
- [dl::Model::test \(C++ function\), 52](#)
- [dl::module::Module \(C++ class\), 45](#)
- [dl::module::Module::~~Module \(C++ function\), 45](#)
- [dl::module::Module::deserialize \(C++ function\), 47](#)
- [dl::module::Module::forward \(C++ function\), 45](#)
- [dl::module::Module::forward\\_args \(C++ function\), 46](#)
- [dl::module::Module::get\\_output\\_shape \(C++ function\), 45](#)
- [dl::module::Module::get\\_outputs\\_index \(C++ function\), 45](#)
- [dl::module::Module::get\\_param\\_memory\\_size \(C++ function\), 46, 47](#)
- [dl::module::Module::inplace \(C++ member\), 47](#)
- [dl::module::Module::m\\_inputs\\_index \(C++ member\), 47](#)
- [dl::module::Module::m\\_outputs\\_index \(C++ member\), 47](#)
- [dl::module::Module::Module \(C++ function\), 45](#)
- [dl::module::Module::name \(C++ member\), 47](#)
- [dl::module::Module::preload \(C++ function\), 46](#)
- [dl::module::Module::print \(C++ function\), 46](#)
- [dl::module::Module::quant\\_type \(C++ member\), 47](#)
- [dl::module::Module::reset \(C++ function\), 46](#)
- [dl::module::Module::retrieve\\_inputs \(C++ function\), 45](#)
- [dl::module::Module::retrieve\\_inputs\\_shape \(C++ function\), 45](#)
- [dl::module::Module::run \(C++ function\), 46](#)
- [dl::module::Module::set\\_preload\\_addr \(C++ function\), 46](#)
- [dl::module::ModuleCreator \(C++ class\), 48](#)
- [dl::module::ModuleCreator::clear \(C++ function\), 48](#)
- [dl::module::ModuleCreator::create \(C++ function\), 48](#)
- [dl::module::ModuleCreator::Creator \(C++ type\), 48](#)
- [dl::module::ModuleCreator::get\\_instance \(C++ function\), 48](#)
- [dl::module::ModuleCreator::print \(C++ function\), 48](#)
- [dl::module::ModuleCreator::register\\_dl\\_modules \(C++ function\), 48](#)
- [dl::module::ModuleCreator::register\\_module \(C++ function\), 48](#)
- [dl::TensorBase \(C++ class\), 39](#)
- [dl::TensorBase::~~TensorBase \(C++ function\), 39](#)
- [dl::TensorBase::assign \(C++ function\), 39, 40](#)
- [dl::TensorBase::auto\\_free \(C++ member\), 44](#)
- [dl::TensorBase::axis\\_offset \(C++ member\), 44](#)
- [dl::TensorBase::cache \(C++ member\), 44](#)
- [dl::TensorBase::caps \(C++ member\), 44](#)
- [dl::TensorBase::compare\\_elements \(C++ function\), 42](#)
- [dl::TensorBase::data \(C++ member\), 44](#)
- [dl::TensorBase::dtype \(C++ member\), 44](#)
- [dl::TensorBase::equal \(C++ function\), 42](#)
- [dl::TensorBase::exponent \(C++ member\), 44](#)

- `dl::TensorBase::flip (C++ function)`, 41  
`dl::TensorBase::get_aligned_bytes (C++ function)`, 40  
`dl::TensorBase::get_aligned_size (C++ function)`, 40  
`dl::TensorBase::get_bytes (C++ function)`, 40  
`dl::TensorBase::get_caps (C++ function)`, 41  
`dl::TensorBase::get_dtype (C++ function)`, 41  
`dl::TensorBase::get_dtype_bytes (C++ function)`, 40  
`dl::TensorBase::get_dtype_string (C++ function)`, 40  
`dl::TensorBase::get_element (C++ function)`, 43  
`dl::TensorBase::get_element_coordinates (C++ function)`, 43  
`dl::TensorBase::get_element_index (C++ function)`, 43  
`dl::TensorBase::get_element_ptr (C++ function)`, 40  
`dl::TensorBase::get_exponent (C++ function)`, 41  
`dl::TensorBase::get_shape (C++ function)`, 40  
`dl::TensorBase::get_size (C++ function)`, 40  
`dl::TensorBase::is_same_shape (C++ function)`, 41  
`dl::TensorBase::pad (C++ function)`, 42  
`dl::TensorBase::preload (C++ function)`, 43  
`dl::TensorBase::print (C++ function)`, 43  
`dl::TensorBase::reset_bias_layout (C++ function)`, 43  
`dl::TensorBase::reshape (C++ function)`, 41  
`dl::TensorBase::set_element_ptr (C++ function)`, 40  
`dl::TensorBase::set_preload_addr (C++ function)`, 43  
`dl::TensorBase::set_shape (C++ function)`, 41  
`dl::TensorBase::shape (C++ member)`, 43  
`dl::TensorBase::size (C++ member)`, 43  
`dl::TensorBase::slice (C++ function)`, 42, 44  
`dl::TensorBase::TensorBase (C++ function)`, 39  
`dl::TensorBase::transpose (C++ function)`, 41  
`DL_LOG_INFER_LATENCY_ARRAY_END (C macro)`, 49  
`DL_LOG_INFER_LATENCY_ARRAY_END_PRINT (C macro)`, 49  
`DL_LOG_INFER_LATENCY_ARRAY_INIT (C macro)`, 49  
`DL_LOG_INFER_LATENCY_ARRAY_INIT_WITH_SIZE (C macro)`, 49  
`DL_LOG_INFER_LATENCY_ARRAY_PRINT (C macro)`, 49  
`DL_LOG_INFER_LATENCY_ARRAY_START (C macro)`, 49  
`DL_LOG_INFER_LATENCY_END (C macro)`, 49  
`DL_LOG_INFER_LATENCY_END_PRINT (C macro)`, 49  
`DL_LOG_INFER_LATENCY_INIT (C macro)`, 49  
`DL_LOG_INFER_LATENCY_INIT_WITH_SIZE (C macro)`, 49  
`DL_LOG_INFER_LATENCY_PRINT (C macro)`, 49  
`DL_LOG_INFER_LATENCY_START (C macro)`, 49
- ## F
- `fbs::FbsLoader (C++ class)`, 59  
`fbs::FbsLoader::~~FbsLoader (C++ function)`, 59  
`fbs::FbsLoader::FbsLoader (C++ function)`, 59  
`fbs::FbsLoader::get_model_num (C++ function)`, 60  
`fbs::FbsLoader::list_models (C++ function)`, 60  
`fbs::FbsLoader::load (C++ function)`, 59, 60  
`fbs::FbsModel (C++ class)`, 60  
`fbs::FbsModel::~~FbsModel (C++ function)`, 61  
`fbs::FbsModel::clear_map (C++ function)`, 65  
`fbs::FbsModel::FbsModel (C++ function)`, 60  
`fbs::FbsModel::get_graph_inputs (C++ function)`, 64  
`fbs::FbsModel::get_graph_outputs (C++ function)`, 64  
`fbs::FbsModel::get_model_doc_string (C++ function)`, 65  
`fbs::FbsModel::get_model_name (C++ function)`, 65  
`fbs::FbsModel::get_model_size (C++ function)`, 65  
`fbs::FbsModel::get_model_version (C++ function)`, 65  
`fbs::FbsModel::get_operation_attribute (C++ function)`, 61, 62  
`fbs::FbsModel::get_operation_inputs_and_outputs (C++ function)`, 63  
`fbs::FbsModel::get_operation_lut (C++ function)`, 63  
`fbs::FbsModel::get_operation_output_shape (C++ function)`, 62  
`fbs::FbsModel::get_operation_parameter (C++ function)`, 63  
`fbs::FbsModel::get_operation_type (C++ function)`, 63  
`fbs::FbsModel::get_tensor_dtype (C++ function)`, 63  
`fbs::FbsModel::get_tensor_exponents (C++ function)`, 63  
`fbs::FbsModel::get_tensor_raw_data (C++ function)`, 63  
`fbs::FbsModel::get_tensor_shape (C++ function)`, 63



`fbs::FbsModel::get_test_input_tensor`  
(C++ *function*), 64

`fbs::FbsModel::get_test_input_tensor_raw_data`  
(C++ *function*), 64

`fbs::FbsModel::get_test_output_tensor`  
(C++ *function*), 64

`fbs::FbsModel::get_test_output_tensor_raw_data`  
(C++ *function*), 64

`fbs::FbsModel::get_test_outputs_name`  
(C++ *function*), 64

`fbs::FbsModel::get_value_info_dtype`  
(C++ *function*), 64

`fbs::FbsModel::get_value_info_exponent`  
(C++ *function*), 64

`fbs::FbsModel::get_value_info_shape`  
(C++ *function*), 64

`fbs::FbsModel::is_parameter` (C++ *function*), 63

`fbs::FbsModel::load_map` (C++ *function*), 65

`fbs::FbsModel::m_param_copy` (C++ *member*), 65

`fbs::FbsModel::memory_addr_in_model`  
(C++ *function*), 65

`fbs::FbsModel::print` (C++ *function*), 61

`fbs::FbsModel::topological_sort` (C++  
*function*), 61