

Tachy-Compile Guidebook

Agenda

1. Overview

- 1-1) What is the Tachy Compile?
- 1-2) Key Features and Advantages

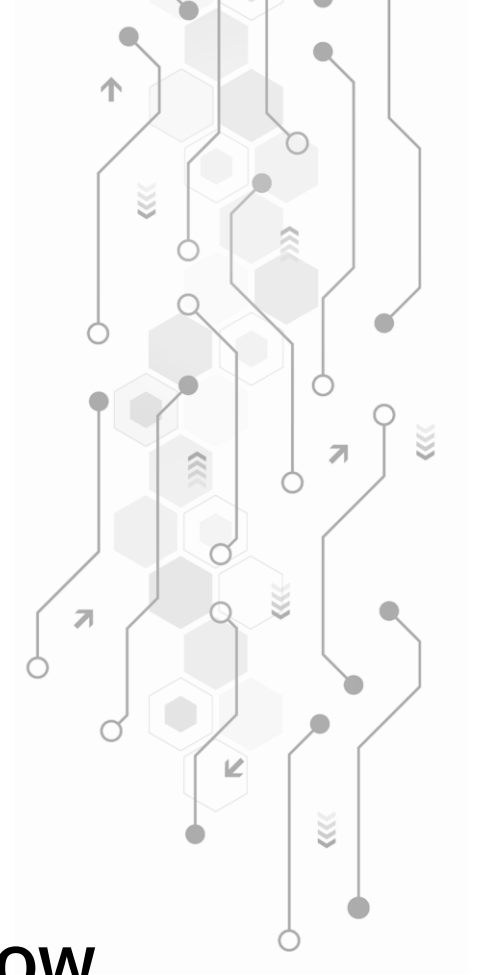
2. Format Optimization Flow

- 2-1) Model Transformation Overview
- 2-2) ONNX conversion process → Tachy format conversion

3. Use Example Guide

〈YOLOv9 Model Learning〉

- 3-1) Preparedness
- 3-2) Main Training
- 3-3) Verifying Learning Results
- 3-4) Optimization Training (optional)
- 3-5) Model Transformation (PyTorch → ONNX → TACHY)

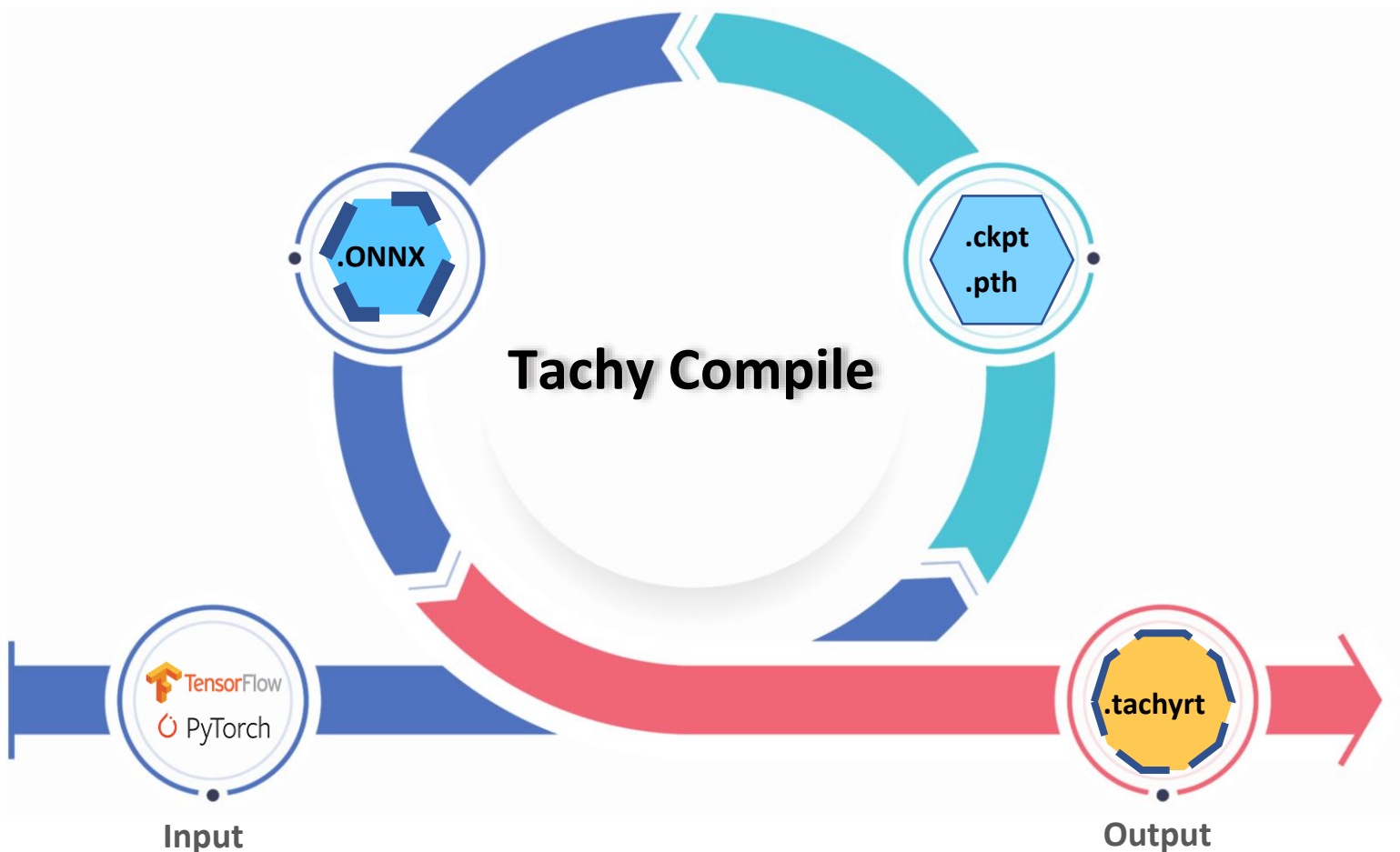


1. Overview

1-1) What is the Tachy Compile?

Tachy Compile is a compiler that transforms models learned from various deep learning frameworks (TensorFlow, PyTorch, etc.) to be optimized and run in high-performance NPU environments. Specifically, by standardizing the model around the Open Neural Network Exchange (ONNX) format and converting it into a TACHY-only format (.tachyrt), it supports lightweight, high speed, and real-time inference.

<Format Optimization Flowchart>



1-2) Key Features and Advantages

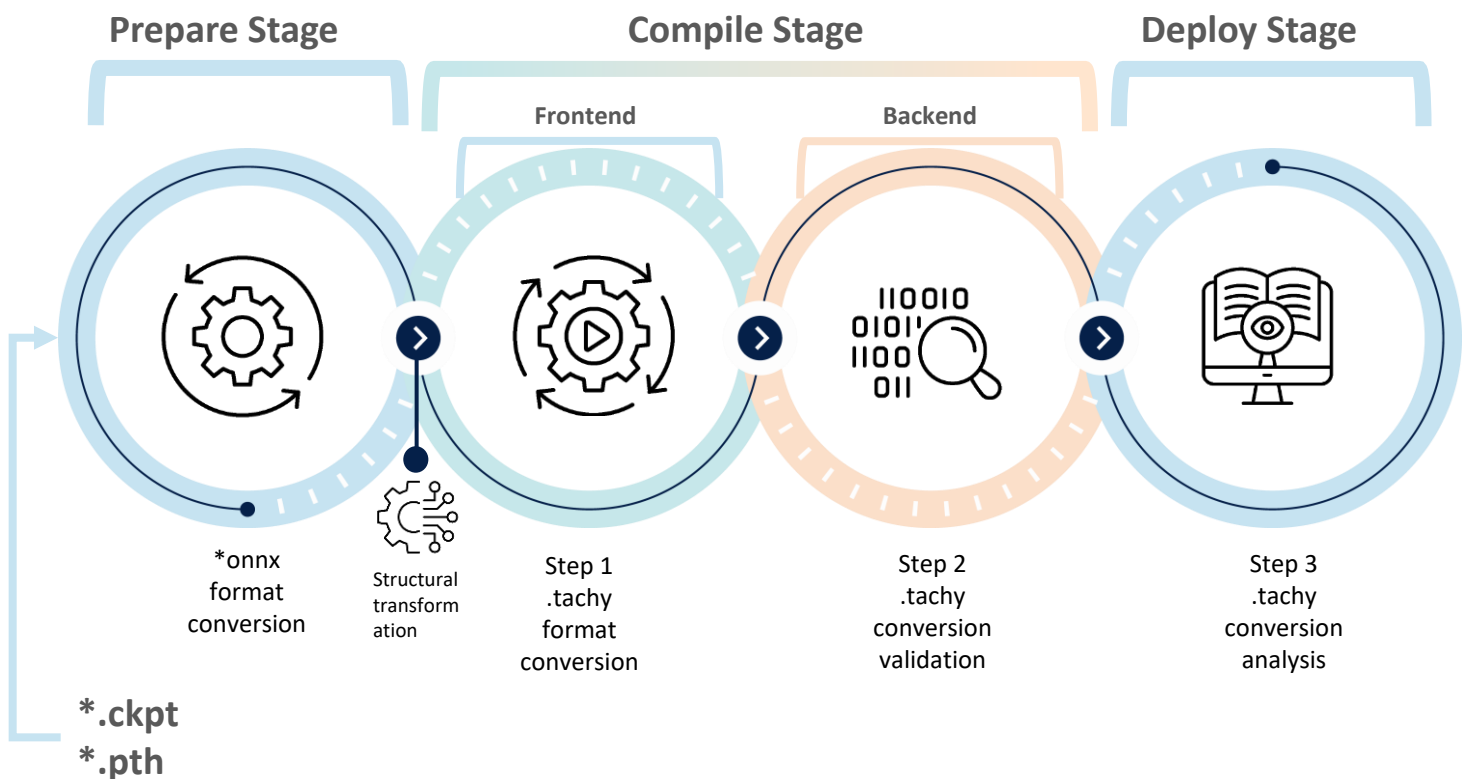
- Consolidate learning, transformation, optimization, and deployment processes into one flow
- Graph optimization and memory weight reduction
- Support and compatibility of a variety of platforms
- Optimization for NPU Application

2. Format Optimization Flow

2-1) Model Transformation Overview

The User prepares a pre-trained model with XWN in TensorFlow or PyTorch, converts it into ONNX format, and finally creates a .tachy file through the Tachy Compiler to optimize the file for application to Deeper-I NPU boards.

2-2) ONNX conversion process → Tachy format conversion



After converting the learned model to ONNX, the ONNX modifier performs operator structure modification and unnecessary node removal. The optimized ONNX model maximizes execution efficiency through XWN lightening. The Tachy Compiler then finally reconstructs the memory structure and operation flow and converts it into .tachy format.

3-1) Preparedness

☐ Building a Learning Environment

Check GPU and CUDA compatibility

```
$ nvidia-smi
```

Check the CUDA version that your GPU supports.

☐ Preparing Dataset

```
{dataset_path}/
├── train/
│   ├── images/
│   └── labels/
├── val/
│   ├── images/
│   └── labels/
└── test/
    ├── images/
    └── labels/
```

☐ Label Format (YOLO format)

```
$ class x_center y_center width height
```

3-2) Main Training

☐ Create Dataset YAML File

```
$ git clone http://asdf.com
$ cd yolov9
```

```
$ cd data
$ vi {DATASET_NAME.yaml} ✖ Create a New File
```

```
path: /mnt/DPI-SHARED/tom/20240625_yolov9/dataset/FP # dataset root dir
train: train/images
val: val/images
test: test/images

# Classes
names:
  0: person
  1: car
  2: bus
  3: truck
```

☐ Create a Model YAML File

```
$ cd ../models/detect
$ cp bsnet-t.yaml {bsnet-t-DATASET_NAME.yaml}
$ vi {bsnet-t-DATASET_NAME.yaml}
```

modification: Correction: Change the nc (number of classes) value to match the dataset

```
# YOLOv9
# parameters
nc: 4 # number of classes
```

Example {bsnet-t-DATASET_NAME.yaml}

❑ Set up a Learning Script

```
$ cd ../../  
$ vi train ✖ Create a New File
```

ex) train command - open the train file, paste the code at the bottom

```
python train.py --workers 8 \  
    --device 0 \  
    --batch 8 \  
    --data data/{DATASET_NAME.yaml} \  
    --img 416 \  
    --cfg models/detect/bsnet-t-{DATASET_NAME.yaml} \  
    --weights " \  
    --name {guide} \  
    --hyp hyp.pretrained.yaml \  
    --min-items 0 \  
    --epochs 300 \  
    --close-mosaic 15
```

❑ Run Learning

```
$ ./train    Execute a file
```

Result: Create running/train/{guide}/weights/best.pt

Epoch	GPU_mem	box_loss	cls_loss	dfl_loss	Instances	Size	LR		
0/299	0.472G	1.543	8.851	0	7	416	0.00874:	100%	15/15 00:04
	Class	Images	Instances	P	R	mAP50	mAP50-95:	100%	2/2 00:00
	all	5	47	0	0	0	0		
Epoch	GPU_mem	box_loss	cls_loss	dfl_loss	Instances	Size	LR		
1/299	0.478G	1.656	8.459	0	38	416	0.007389:	100%	15/15 00:01
	Class	Images	Instances	P	R	mAP50	mAP50-95:	100%	2/2 00:00
	all	5	47	0	0	0	0		
Epoch	GPU_mem	box_loss	cls_loss	dfl_loss	Instances	Size	LR		
2/299	0.478G	1.501	8.652	0	49	416	0.006037:	100%	15/15 00:01
	Class	Images	Instances	P	R	mAP50	mAP50-95:	100%	2/2 00:00
	all	5	47	0.00405	0.0213	0.00207	0.00165		
Epoch	GPU_mem	box_loss	cls_loss	dfl_loss	Instances	Size	LR		
3/299	0.478G	1.5	8.168	0	20	416	0.004684:	100%	15/15 00:01
	Class	Images	Instances	P	R	mAP50	mAP50-95:	100%	2/2 00:00
	all	5	47	0.0287	0.149	0.0165	0.0135		

3-3) Verifying Learning Results

❑ Setting up and running verification scripts

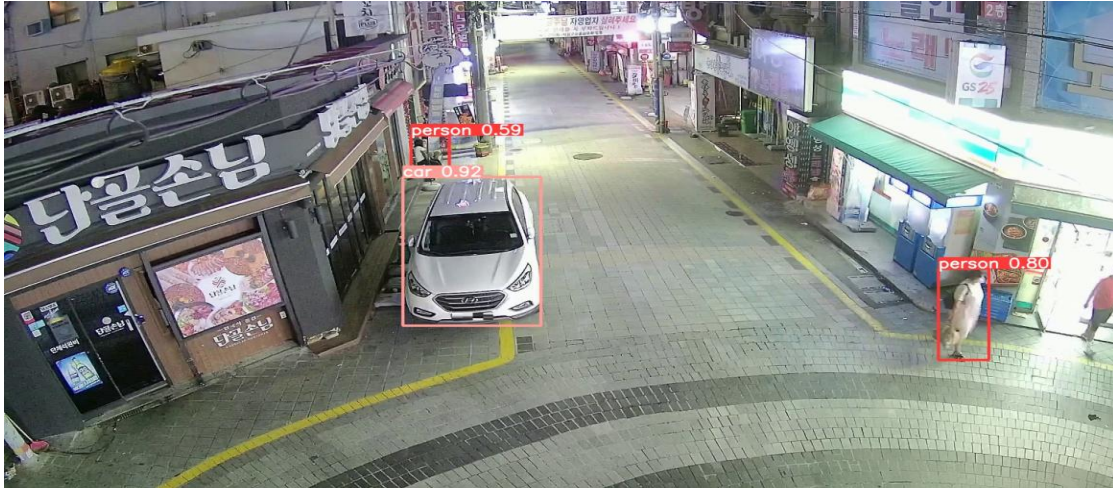
```
$ vi test ✖ Create a New File
```

ex) test command - open test file, paste code at the bottom

```
python detect.py --source {dataset_path}/test/images \  
    --img 416 \  
    --device 0 \  
    --weights runs/train/{guide}/weights/best.pt \  
    --name {guide} \  
    --save-txt
```

```
$ ./test    Execute a file
```

Results: Save detection results in the runs/detect/{guide} directory



Example detect image

3-4) Optimization Training (optional)

(Optional: When conducting quantization/optimization training)

☐ Create Model YAML for Optimization

```
$ cd ../models/detect
$ cp bsnet-t-o.yaml {bsnet-t-o-DATASET_NAME.yaml}
$ vi {bsnet-t-o-DATASET_NAME.yaml}
```

```
# YOLOv9

# parameters
nc: 4 # number of classes
```

Example {bsnet-t-o-DATASET_NAME.yaml}

```
11 anchors: 3
12
13 # getnn backbone
14 backbone:
15 [
16   # conv down
17   [1, 1, conv, [16, 1, 1], [4, null], # 0-P3/2
18   [1, 1, conv, [16, 1, 1], [4, null], # 1
19   [1, 1, conv, [16, 1, 1], [4, null], # 2
20   [1, 1, conv, [16, 1, 1], [4, null], # 3
21   [1, 1, conv, [16, 1, 1], [4, null], # 4-P3/4
22   [1, 1, conv, [16, 1, 1], [4, null], # 5
23   [1, 1, conv, [16, 1, 1], [4, null], # 6
24   [1, 1, conv, [16, 1, 1], [4, null], # 7
25   [1, 1, conv, [16, 1, 1], [4, null], # 8-P3/8
26   [1, 1, conv, [16, 1, 1], [4, null], # 9
27   [1, 1, conv, [16, 1, 1], [4, null], # 10
28   [1, 1, conv, [16, 1, 1], [4, null], # 11
29   [1, 1, conv, [16, 1, 1], [4, null], # 12-P4/16
30   [1, 1, conv, [16, 1, 1], [4, null], # 13
31   [1, 1, conv, [16, 1, 1], [4, null], # 14
32   [1, 1, conv, [16, 1, 1], [4, null], # 15
33   [1, 1, conv, [16, 1, 1], [4, null], # 16
34   [1, 1, conv, [16, 1, 1], [4, null], # 17
35   [1, 1, conv, [16, 1, 1], [4, null], # 18-P5/32
36   [1, 1, conv, [16, 1, 1], [4, null], # 19-P5/32
37   [1, 1, conv, [16, 1, 1], [4, null], # 20
38 ]
39
40 head:
41 [
42   # P5
43   [1, 1, conv, [16, 1, 1], [4, null], # 21
44   [1, 1, conv, [16, 1, 1], [4, null], # 22
45 ]
46 # P4
47 [1, 1, conv, [16, 1, 1], [4, null], # 23
48 [1, 1, conv, [16, 1, 1], [4, null], # 24
49 [1, 1, conv, [16, 1, 1], [4, null], # 25
50 [1, 1, conv, [16, 1, 1], [4, null], # 26
51 ]
52 # P3
53 [1, 1, conv, [16, 1, 1], [4, null], # 27
54 [1, 1, conv, [16, 1, 1], [4, null], # 28
55 [1, 1, conv, [16, 1, 1], [4, null], # 29
56 [1, 1, conv, [16, 1, 1], [4, null], # 30
57 [1, 1, conv, [16, 1, 1], [4, null], # 31-P3/16
58 [1, 1, conv, [16, 1, 1], [4, null], # 32-P3/16
59 ]
```

diff main training model.yaml, opt training model.yaml

opt training differ bit : int / 4 / Quantization range (bit-1 ** 2)

□ Setting Up Optimization Learning Scripts

```
$ cd ../../  
$ vi train ✖ Create a New File
```

ex) train command - open the train file, paste the code at the bottom

```
python train.py --workers 8 \  
    --device 0 \  
    --batch 8 \  
    --data data/ {DATASET_NAME.yaml} \  
    --img 416 \  
    --cfg models/detect/ {bsnet-t-o-DATASET_NAME.yaml} \  
    --weights 'runs/train/{guide}/weights/best.pt' \  
    --name {guide-opt} \  
    --hyp hyp.optimization.yaml \  
    --min-items 0 \  
    --epochs 100 \  
    --close-mosaic 15
```

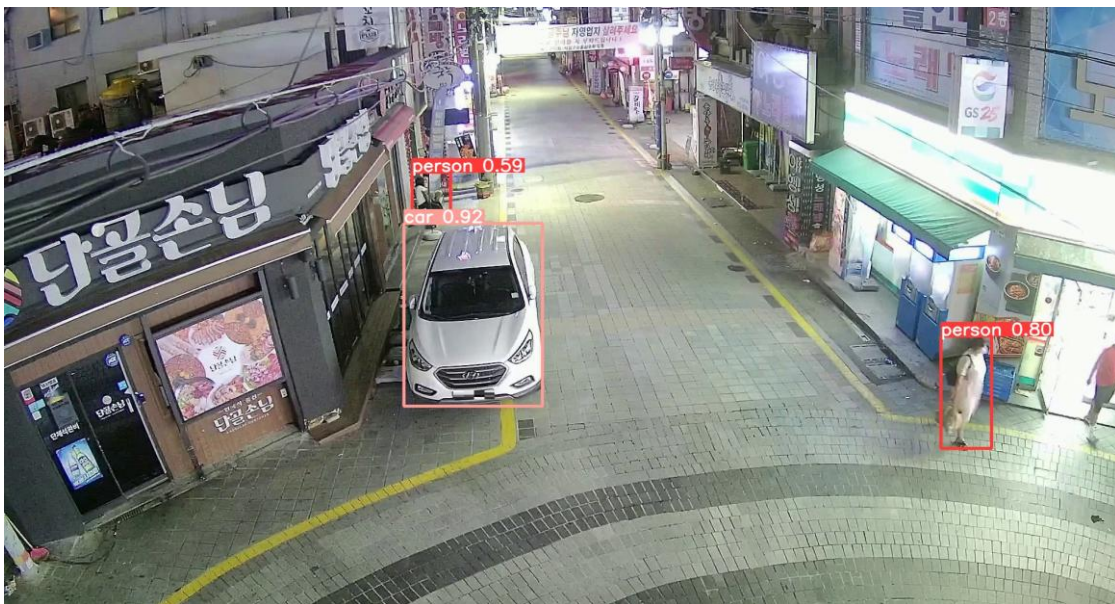
* Pay attention to yaml file name in --data and --cfg options

□ Optimization Learning and Validation

```
$ ./train Execute a file
```

□ Reverify optimized results

```
$ vi test  
  
python detect.py --source {/dataset_path/test/images} \  
    --img 416 --device 0 \  
    --weights 'runs/train/{guide-opt}/weights/best.pt' \  
    --name {guide-opt} --save-txt  
  
$ ./test Execute a file
```



3-5) Model Transformation(PyTorch → ONNX → TACHY)

What is Compile?

A batch script that automatically compiles the yolov9 model learned by pt, etc. into a tachy runtime package (.tachyrt, param.bin, etc.) that tachy-bs npu can understand.

☐ Modify the compile script

Modify the compile script inside the yolov9 project.

vi compile ✕ Open a file

<Explanation of key variables in the compile script> ✕ Required Settings

config	default	explanation
Input path		
PRE_PARAM_DIR	runs/golden/drone/bsnet-t-drone3c/weights	Directory path where the pre-optimization model weight (.pt) is located
OPT_PARAM_DIR	runs/golden/drone/bsnet-t-o-drone3c/weights	Directory path where model weight (.pt) is located after optimization
Size		
ONNX_INPUT_SHAPE	"1 3 320 416"	Input tensor form during ONNX conversion (B, C, H, W)
TACHY_INPUT_SHAPE	"1 320 416 3"	Input tensor form for TACHY-Compiler (B, H, W, C)

< Option Settings >

config	default	explanation
TRAIN_FILE_TYPE	"pt_yolov9"	Learning File Format for Deployment Scripts (PyTorch-YOLOv9)
ONNX_FILE_TYPE	"pt2onnx"	PyTorch → ONNX Transformation Script Type
SRC_PRE_PATH	os.path.join(PRE_PARAM_DIR, "best.pt")	Full path to weight files before optimization
SRC_OPT_PATH	os.path.join(OPT_PARAM_DIR, "best.pt")	Full path to weighted files after optimization
WORK_DIR	datetime.datetime.now().strftime("20%ym%d_%H%M%S")	The name of the working directory that is generated at the time of execution (for example, 20250429_153240)
TEMP_DIR	os.path.join(WORK_DIR, "temp")	Temporary file storage folder inside the working directory
DST_EVAL_PATH	os.path.join(WORK_DIR, "eval.pt")	PyTorch Weight Storage Path for Evaluation (Uncompressed)
DST_COMP_PATH	os.path.join(WORK_DIR, "comp.pt")	Compressed (compiled) PyTorch weight storage path
DST_TEMP_PATH	"" (Empty String)	Paths filled when temporary deliverables are needed (currently unused)
ONNX_EVAL_PATH	os.path.join(WORK_DIR, "eval.onnx")	EVAL ONNX model storage path
ONNX_COMP_PATH	os.path.join(WORK_DIR, "comp.onnx")	COMP ONNX model storage path
LAYER_EVAL_PATH	os.path.join(WORK_DIR, "eval.tachy")	Layer-level EVALTACHY file path
LAYER_COMP_PATH	os.path.join(WORK_DIR, "comp.tachy")	Layer-level COMP TACHY file path
L_DEFAULT_PAD_MODE	--default_pad_mode=dynamic	Padding Mode Setting Flag at Layer Compilation
BLOCK_COMP_PATH	os.path.join(WORK_DIR, "block.tachy")	Block-level compression TACHY file path
B_DEFAULT_PAD_MODE	--default_pad_mode=dynamic	Padding mode setting flag when compiling blocks
DEFAULT_PAD_ORDER	--default_pad_order=up_left	Specify the default direction to apply padding
LOGIT_ORDER	--logit_order 5 3 0 4 2 1	Output logit (class) order rearrangement options

☐ **Execute a file**

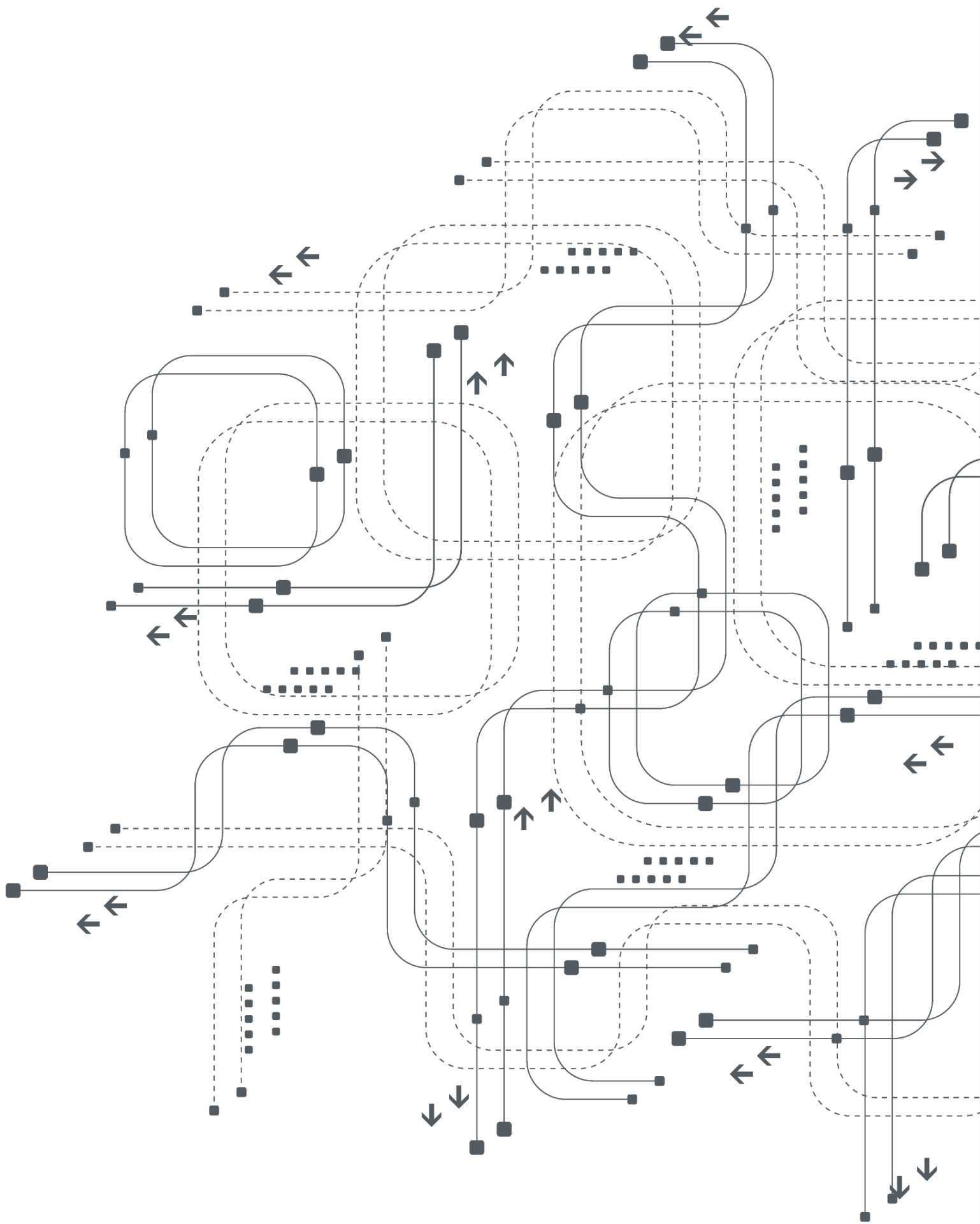
```
./compile
```

☐ **Check the results**

```
ls 20250429_152340/*
```

☐ **Example Output**

```
model_224x224x3_inv-t.tachyrt inst.bin param.bin inst.json
```



Customer Support

- Business : partner@deeper-i.ai
- website : <https://www.deeper-i.ai>