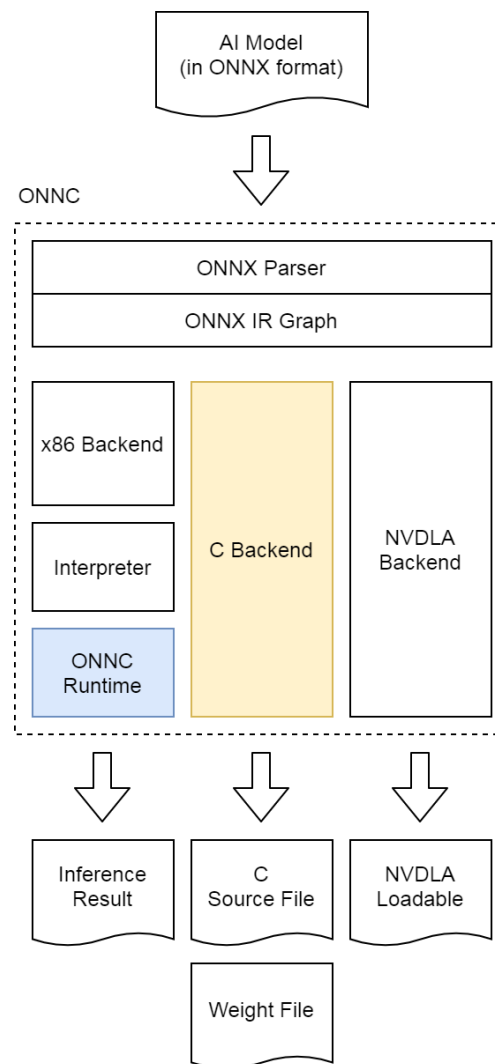# **Skymizer** | Introduction of ONNC C Backend

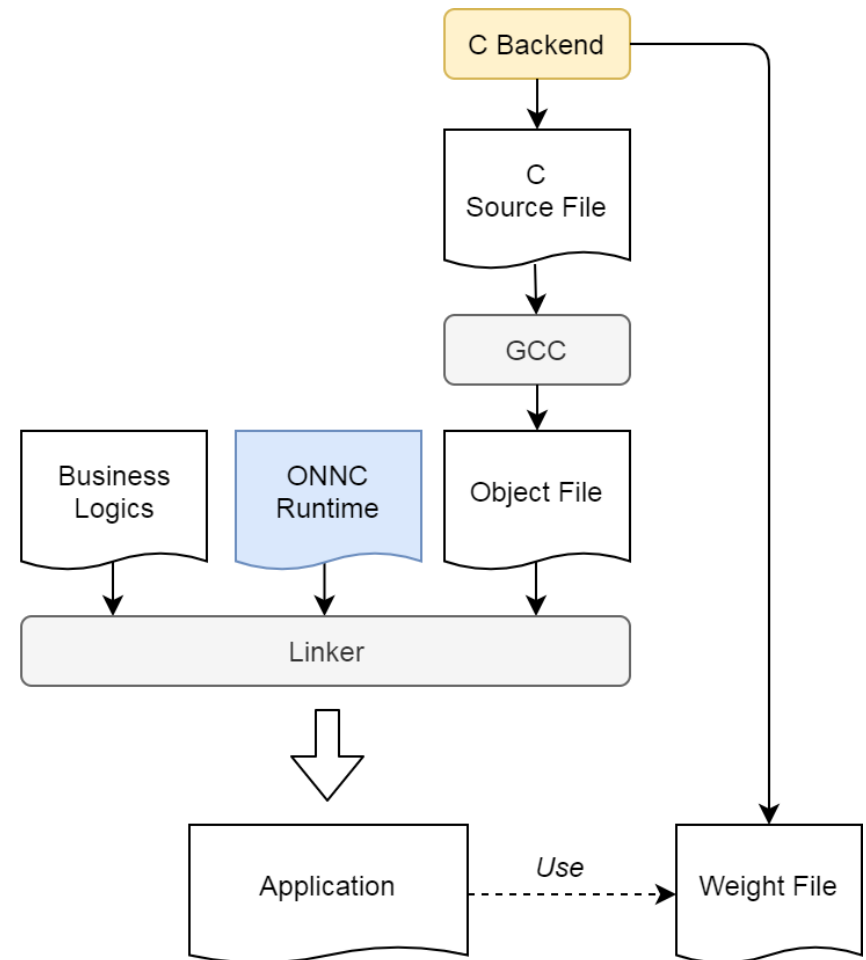Po-Yen Chen (poyenc@skymizer.com)

a

# ONNC Framework

- ONNC Runtime is a model-inference library

- C Backend transforms a model into a C source file and a weight file

- Users can run model inference by using several of C functions

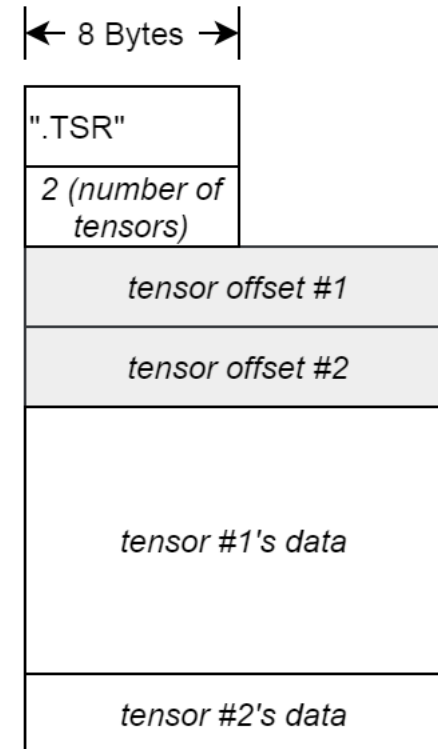**skymizer**

a

# C Backend Workflow

- C Backend creates implementation of a single function: *model_main()*

- *model_main()* is the entry function for input model inference flow

- *model_main()* employ ONNC Runtime to compute results

**skymizer**

a

# Weight file format

```
12   struct ONNC_RUNTIME_tensor_offset {
13     uint64_t offset; /* Tensor offset */
14     uint64_t size;   /* Size of tensor in bytes */
15   };
16
17   #ifndef ONNC_RUNTIME_TENSOR_FILE_MAGIC
18   #   define ONNC_RUNTIME_TENSOR_FILE_MAGIC ".TSR"
19   #endif
20
21   struct ONNC_RUNTIME_tensor_offset_table
22   {
23     char                        magic[8]; /* Tensor File magic number. */
24     uint64_t                    number_of_tensors;
25     struct ONNC_RUNTIME_tensor_offset tensor_offsets[];
26   };
27
```
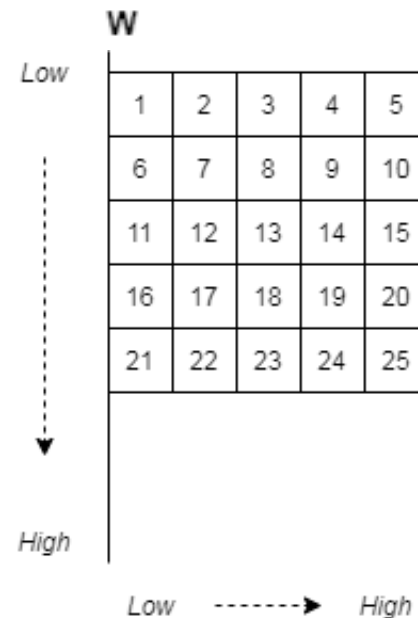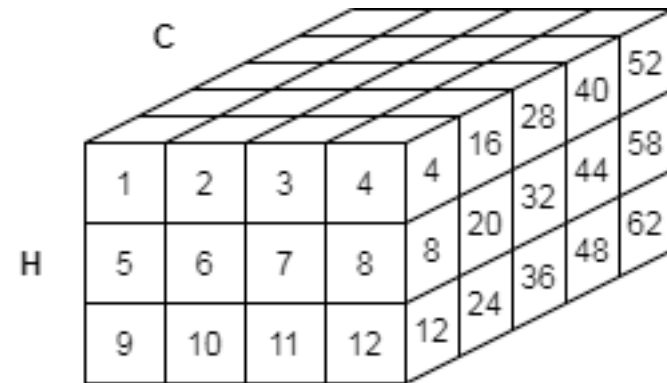


← 8 Bytes →

".TSR"

2 (number of tensors)

tensor offset #1

tensor offset #2

tensor #1's data

tensor #2's data

**skymizer**

a

# Tensor data format

include/onnc/Runtime/onnc-runtime.h:7

```
7    struct ONNC_RUNTIME_tensor_file
8    {
9      void* data; /* Implementation defined data */
10   };
```

include/onnc/Runtime/onnc-runtime.h:28

```
28   struct ONNC_RUNTIME_tensor_view
29   {
30     void*    data;
31     uint64_t size; /* Size of tensor in bytes */
32   };
```

**skymizer**

a

# Tensor access API (user defined)

include/onnc/Runtime/onnc-runtime.h:39

```
39   // Client Library
40   const struct ONNC_RUNTIME_tensor_offset_table*
41                           ONNC_RUNTIME_read_tensor_offset_table(struct ONNC_RUNTIME_tensor_file* file);
42   struct ONNC_RUNTIME_tensor_view ONNC_RUNTIME_read_tensor(struct ONNC_RUNTIME_tensor_file* file, uint64_t tensor);
```

## A possible implementation

example/runtime/src/client-lib.c

```
16   struct ONNC_RUNTIME_tensor_view ONNC_RUNTIME_read_tensor(struct ONNC_RUNTIME_tensor_file* file, uint64_t tensor)
17   {
18     if (file == NULL) {
19       const struct ONNC_RUNTIME_tensor_view tensor_view = {.data = NULL, .size = 0};
20       return tensor_view;
21     }
22
23     const struct ONNC_RUNTIME_tensor_offset_table* const table = ONNC_RUNTIME_read_tensor_offset_table(file);
24     if (!ONNC_RUNTIME_has_tensor(table, tensor)) {
25       const struct ONNC_RUNTIME_tensor_view tensor_view = {.data = NULL, .size = 0};
26       return tensor_view;
27     }
28
29     const struct ONNC_RUNTIME_tensor_offset tensor_offset = ONNC_RUNTIME_get_tensor_offset(table, tensor);
30     const struct ONNC_RUNTIME_tensor_view   tensor_view   = {.data = (char*)file->data + tensor_offset.offset,
31                                                              .size = tensor_offset.size};
32     return tensor_view;
33   }
```

**skymizer**

a

# Context and model_main()

```
45   struct ONNC_RUNTIME_inference_context
46   {
47     struct ONNC_RUNTIME_tensor_file* input;
48     struct ONNC_RUNTIME_tensor_file* weight;
49     uint64_t                         id;
50     void (*completed)(uint64_t id, struct ONNC_RUNTIME_tensor_view output);
51   };
52
53   /**
54    * ONNC generated entry point.
55    * @param context The ONNC Runtime Context.
56    */
57   int model_main(const struct ONNC_RUNTIME_inference_context* context);
```

**skymizer**

a

# User Application Example

example/runtime/src/client-app.c

```c
40   void finish(uint64_t id, struct ONNC_RUNTIME_tensor_view output)
41   {
42     const float* const values = output.data;
43     const size_t       count  = output.size / sizeof(float);
44     printf("[");
45     for (size_t idx = 0; idx < count; ++idx) {
46       printf("%f, ", values[idx]);
47     }
48     printf("]");
49   }
50
51   int main(int argc, char* argv[])
52   {
53     if (argc < 3) {
54       fprintf(stderr, "usage: %s foo.input foo.weight\n", argv[0]);
55       return EXIT_FAILURE;
56     }
57
58     struct ONNC_RUNTIME_tensor_file* const input  = open_tensor_file(argv[1]);
59     struct ONNC_RUNTIME_tensor_file* const weight = open_tensor_file(argv[2]);
60
61     struct ONNC_RUNTIME_inference_context context = {.input = input, .weight = weight, .id = 0, .completed = finish};
62
63     model_main(&context);
64
65     close_tensor_file(input);
66     close_tensor_file(weight);
67
68     return EXIT_SUCCESS;
69   }
```

8

**‹ymizer**

a

# How to build a user application

(Inside onnc-community docker container)

## 1. Prepare ONNC Runtime library and onnx tensor convertor **pb2t**

```
$ sudo cp /onnc/onnc-umbrella/build-normal/lib/Runtime/libonnc-rt.a /usr/local/lib
# go to build directory and sync source files
$ cd /onnc/onnc-umbrella/build-normal && ssync
# build convertor and install it into PATH
$ make pb2t && make install
```

## 2. Build example application and run model inference
example/runtime

```
# prepare build directory
$ cd /onnc/onnc/example/runtime
$ mkdir build && cd build
# prepare service library
$ onnc -mquadruple clang /models/bvlc_alexnet/model.onnx -o ./test.c
# prepare sample input tensor file
$ pb2t /models/bvlc_alexnet/test_data_set_0/input_0.pb ./test.input
$ cp ./test.c ../src/onnc-runtime-service.c
# configure & build example project
$ cmake .. && make
$ ./example/inference test.input test.weight
[0.000043, 0.000046, 0.000024, 0.000011, 0.000114, 0.000469, …
…
…, 0.000035, 0.000148, 0.000964, 0.000134, 0.001431, 0.000448, ]
```

**skymizer**

a

# References

- ONNC C-backend Tutorial

https://github.com/ONNC/onnc/blob/master/docs/ONNC-C-Backend-Guide.md

**skymizer**

# Skymizer Taiwan Inc.

**skymizer**

**skymizer**

Boost deep learning accelerator
with compiler technology

**CONTACT US**

**E-mail**  sales@skymizer.com      **Tel**   +886 2 8797 8337

**HQ**  12F-2, No.408, Ruiguang Rd., Neihu Dist., Taipei City 11492, Taiwan

**BR**   Center of Innovative Incubator, National Tsing Hua University,
Hsinchu Taiwan

https://skymizer.com