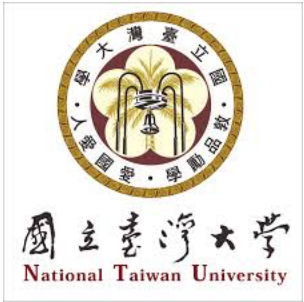# ONNC–CIM | An End-to-End Crossbar-Based CIM  Simulator

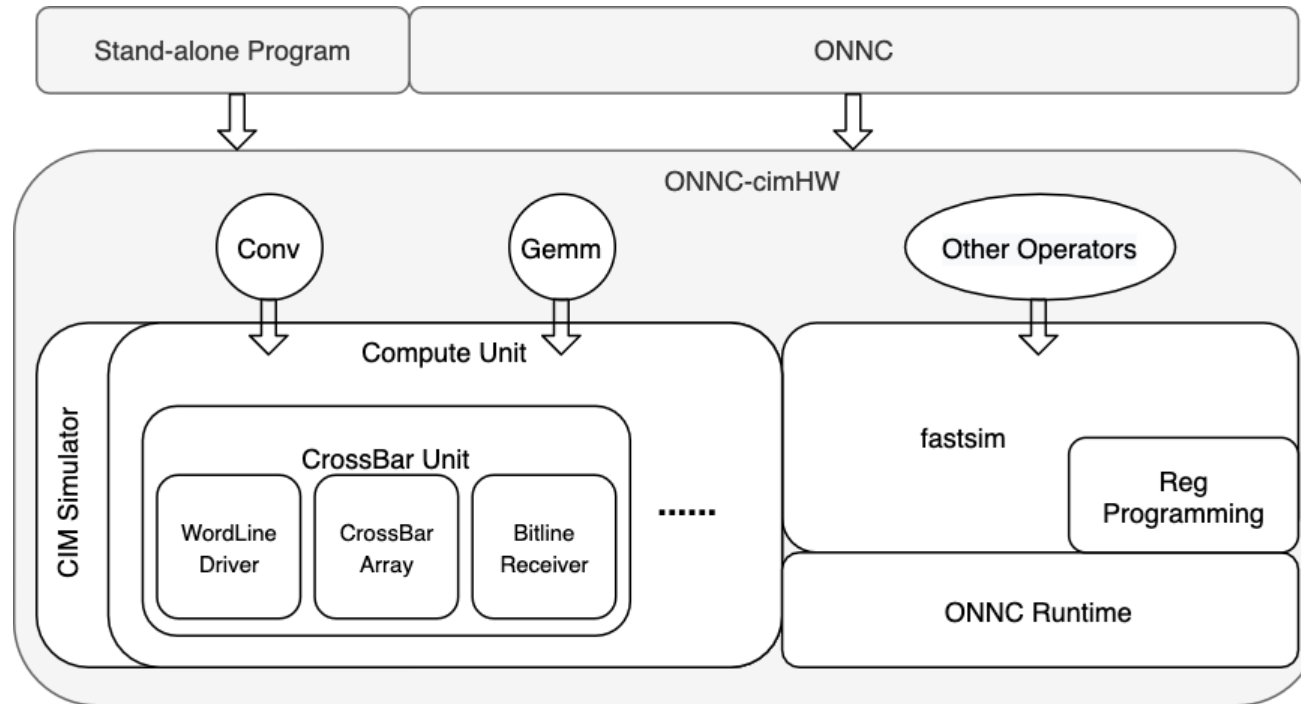Jen-ho Kuo & Yu-Jie Wang l May 2020

Skymizer Taiwan Inc.

# A Joint Work with the Embedded Computing Lab @ CSIE, National Taiwan University

Meng-Yao Lin, Wei-Ting Lin, Tzu-Hsien Yang, I-Ching Tseng, Hsiang-Yun Cheng, Chia-Lin Yang

**NTU CSIE Embedded Computing Lab**
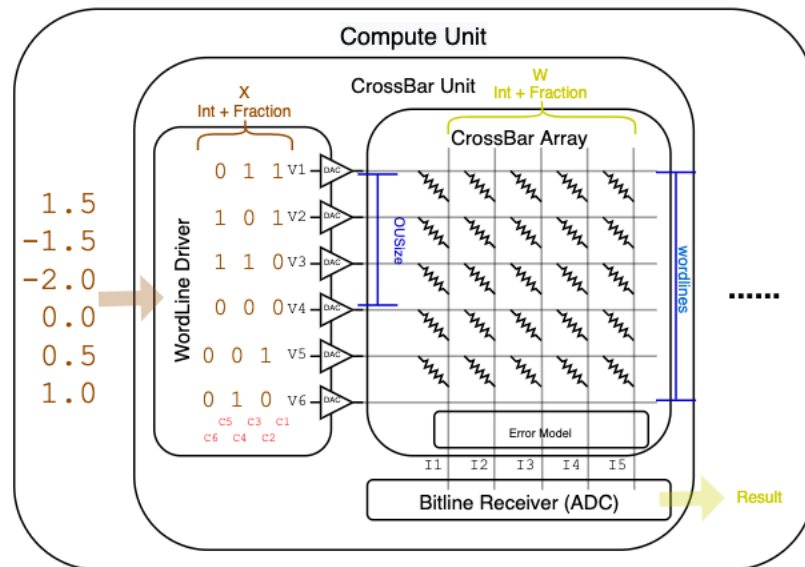
ONNC

# Outline

- Project overview

- Device Error Model

- How to build and run ONNC-CIM

- How to create a user application in the stand-alone mode

- How to support a new operator in ONNC-CIM

ONNC

# Project Overview



- Support ONNX Operator v1.3.0
- Support full-model simulation
  – CrossBar device model for **Conv** and **Gemm**
  – Fastsim (behavior model ) for **Other operators**
- Usages
  – Build as a library in the **Stand-alone** mode
  – Build with an ONNC frontend in the **ONNC** mode
- Extendibility
  – Support adding new operators that use crossbar devices
  – A unified JSON format for device model

# CrossBar Unit



- **CIM simulator**
  - Compute Unit has several Crossbar Units
  - Crossbar Unit has
    - WordLine Driver (DAC)
    - Crossbar Array
    - Bitline Receiver (ADC)

- **CrossBar Array**
  - The resistor could be configured by vector W with given quantized parameters

- **WordLine Driver**
  - Repeatedly accepts vector X and quantizes into given length of bits

- **Bitline Receiver**
  - For each cycle, wordline driver electrifies maximum OUSize of wordlines of bits into crossbar array
  - Amplification will be sensed with error model simulation
  - Final result is then added and shifted cycle-by-cycle

# How to Build and Run ONNC-CIM

- ## Prerequisite
  - Docker
  - ONNC Community Docker image

- ## Clone the onnc-cim project from GitHub
  `$ git clone https://github.com/onnc/onnc-cim.git && cd onnc-cim`

- ## Run the container with
  `$ docker run -it -v `pwd`:/onnc/onnc --cap-add=SYS_PTRACE --security-opt seccomp=unconfined onnc/onnc-community:latest`

# Step to Build  onnc-cim in Stand-alone mode

1. Create build folder and compile onnc-cim
   $ mkdir -p /onnc/onnc-cimHW.build

   $ cd /onnc/onnc-cimHW.build
   $ cmake -DCMAKE_BUILD_TYPE=Release /onnc/onnc/skysim/onnc-cimHW
   $ make -j8

2. Install the built library
   $ sudo make install
   $ sudo ldconfig

ONNC

# Build Options

- DCMAKE_BUILD_TYPE
  - DCMAKE_BUILD_TYPE=Release
  - DCMAKE_BUILD_TYPE=Debug

- DBUILD_TESTS=ON
  - Test cases will be built to a single test_all program in tests folder.
  - ./tests/test_all share/cimConfig.json

ONNC

# Testing onnc-cimHW in stand-alone mode

- You may find a test program, "main", in the example folder
- Run simulation with a device configuration file, "cimConfig.json"
  $ ./example/main ./share/cimConfig.json

```
onnc@2cc187bff70a:/onnc/onnc-cimHW.build$ ./example/main ./share/cimConfig.json
Result:
673.5, 673.5, 673.5, 673.5, 673.5,
```

ONNC

# Step to Build onnc-cim in onnc mode

1. build onnc-cim and install the library in the docker container (page 10)

2. build onnc with the onnc-cim library

```
$ cd /onnc/onnc-umbrella/build-normal
$ smake -j8 install
```
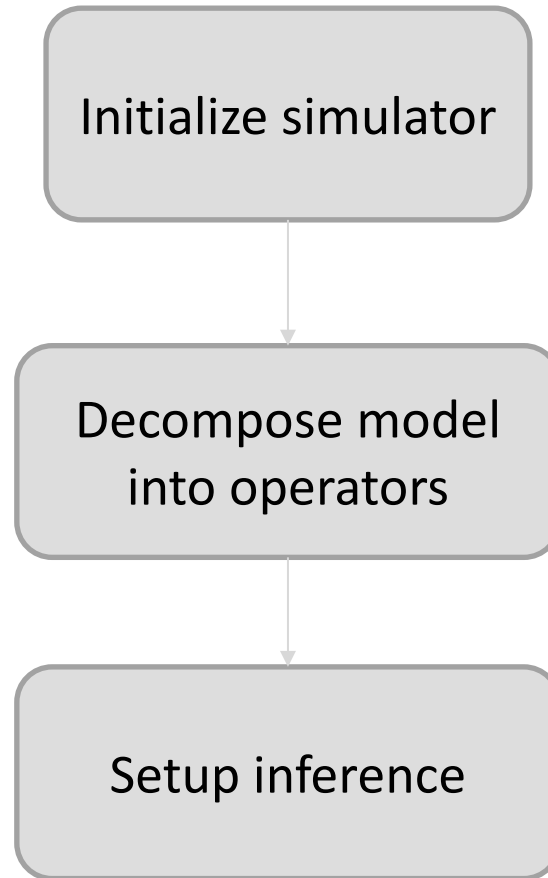
ONNC

# Testing onnc-cim in onnc mode

Example to use nn on MNIST

$ onnc-cim -mquadruple cim /onnc/onnc-cimHW.build/share/nn_mnist/model.onnx --cim-interpreter-input /onnc/onnc-cimHW.build/share/nn_mnist/test_data_set_0/input_2.pb --cim-config /onnc/onnc-cimHW.build/share/cimConfig.json

```
onnc@2cc187bff70a:/onnc/onnc-umbrella/build-normal$ onnc-cim -mquadruple cim /onnc/onnc-cimHW.build/share/nn_mnist/model.onnx --cim-interpreter-input
/onnc/onnc-cimHW.build/share/nn_mnist/test_data_set_0/input_2.pb --cim-config /onnc/onnc-cimHW.build/share/cimConfig.json
[v1] weight memory: 23992
[v1] internal memory: 77312
[v1] total inference time: 954253900 ns
[-42.645618, -225.843338, 1401.953369, -45.595074, -79.142662, -206.296021, -237.601776, -79.743721, -108.739471, -234.470078, ]
```

ONNC

# How to write a program for using the stand-alone library

Initialize simulator

Decompose model into operators

Setup inference

ONNC

# Initialize simulator

- /onnc/onnc/skysim/onnc-cimHW/example/main.cc

# Decompose a model into operators

Input          1x1x5x5

Conv           (1x3x5x5)
               kernel=3x1x3x3, bias=1x3, pads=[1,1,1,1], strides=[1,1])

Reshape        (1x75)
               shape=[1,75]

Gemm           (1x5)
               gemmB=5x75, transA=0, transB=1

Output

ONNC

# Inference Setup

1. Prepare input, output, and activation tensors

2. Specify operator attributes

3. Call onnc-cimHW runtime API and get output tensors

# Inference Setup

Prepare input, output, and activation tensors

▼  📄 **skysim/cimHWSim/example/main.cc** 0 → 100644  📋

```
38  +   // input(1x1x5x5)
39  +   const std::vector<dim_type> input_shape {1, 1, 5, 5};
40  +   element_type intput[] = {
41  +       1.5, 1.0, 0.5, 1.5, 1.0,
42  +       1.5, 1.0, 0.5, 1.5, 1.0,
43  +       1.5, 1.0, 0.5, 1.5, 1.0,
44  +       1.5, 1.0, 0.5, 1.5, 1.0,
45  +       1.5, 1.0, 0.5, 1.5, 1.0
46  +   };
```

ONNC

# Inference Setup

Prepare input, output, and activation tensors

▾ 📄 **skysim/cimHWSim/example/main.cc** 0 → 100644 📋

```
51  +    // kernel(1x3)
52  +    const std::vector<dim_type> weight_shape {3, 1, 3, 3};
53  +    const_element_type weight[] = {
54  +       1.5, 1.0, 0.5,
55  +       1.5, 1.0, 0.5,
56  +       1.5, 1.0, 0.5,
57  +
58  +       1.5, 1.0, 0.5,
59  +       1.5, 1.0, 0.5,
60  +       1.5, 1.0, 0.5,
61  +
62  +       1.5, 1.0, 0.5,
63  +       1.5, 1.0, 0.5,
64  +       1.5, 1.0, 0.5
65  +    };
```

▾ 📄 **skysim/cimHWSim/example/main.cc** 0 → 100644 📋

```
71  +    // conv1(1x3x5x5)
72  +    const std::vector<dim_type> conv1_shape {1, 3, 5, 5};
73  +    element_type conv1[75];
```

**NTU CSIE Embedded Computing Lab**

ONNC

# Inference Setup

Specify operator attributes

```
74  +  // dilation(1x2)
75  +  std::vector<dim_type> dilations {1, 1};
76  +  // group(1)
77  +  const_dim_type group = 1;
78  +  // kernel(1x2)
79  +  std::vector<dim_type> kernel {3, 3};
80  +  // pads(1x4)
81  +  std::vector<dim_type> pads {1, 1, 1, 1};
82  +  // strides(1x2)
83  +  std::vector<dim_type> strides {1, 1};
```

skysim/cimHWSim/example/main.cc 0 → 100644

ONNC

# Inference Setup

Call onnc-cimHW function and get output



```
    skysim/cimHWSim/example/main.cc 0 → 100644

 84  +    ops.cimHW_Conv_float(
 85  +      reinterpret_cast<void*>(&configFile),
 86  +                intput, input_shape.size()  , input_shape.data(),
 87  +                weight, weight_shape.size() , weight_shape.data(),
 88  +                bias  , bias_shape.size()   , bias_shape.data(),
 89  +                conv1 , conv1_shape.size()  ,  conv1_shape.data(),
 90  +      "",
 91  +      dilations.data(), dilations.size(),
 92  +      group,
 93  +      kernel.data(), kernel.size(),
 94  +      pads.data(), pads.size(),
 95  +      strides.data(), strides.size()
 96  +    );
```

ONNC

# How to implement a new operator

1. Create a new operator class from the template
2. Register to desired factory class
3. Add the new file to the building system
4. Rebuild

ONNC

# Create a new operator class from the template

- $ cp lib/op.template.hh lib/matmul.hh

- replace ${OperatorName} => MatMul

- replace ${argument_list}
  - skysim/onnc-cimHW/include/abstractOpFactory.hh
    - 2<sup>nd</sup> parameter of create${OperatorName}Op function

```
,const_element_type* input_A
,const_dim_type input_A_ndim, const_dim_type* input_A_dims
,const_element_type* input_B
,const_dim_type input_B_ndim, const_dim_type* input_B_dims
,element_type* output_Y
,const_dim_type output_Y_ndim, const_dim_type* output_Y_dims
```

```
 8    + #pragma once
 9    + #include "hardware/ComputeUnit.hh"
10    + #include <cimHWOp.hh>
11    +
12    + namespace cimHW {
13    +
14    + class cimHW${OperatorName}Op : public cimHWOp {
15    + public:
16    +   cimHW${OperatorName}Op(
17    +     const std::string &configFile
18    +     ${argument_list}
19    +   );
20    +
21    +   virtual void simulate() final;
22    +
23    + private:
24    +
25    +   // Use a compute unit
26    +   ComputeUnit              m_cimCU;
27    + };
28    +
29    + } // namespace cimHW
        \ No newline at end of file
```

ONNC

# Create a new operator class from the template

- op.template.hh
  - Add member for computing & output

```
19  +  );
20  +
21  +    virtual void simulate() final;
22  +
23  + private:
24  +
25  +    // Use a compute unit
26  +    ComputeUnit                m_cimCU;
27  + };
28  +
29  + } // namespace cimHW
      \ No newline at end of file
```

```
element_type*          m_output_Y;

MatrixXfRowMajor       m_matrix_A;
MatrixXfRowMajor       m_matrix_B;
```

ONNC

# Create a new operator class from the template

- $ cp lib/op.template.cc lib/matmul.cc
- replace ${OperatorName}
- replace ${argument_list}
  - skysim/onnc-cimHW/include/abstractOpFactory.hh
    - 2nd parameter of create${OperatorName}Op function

```
,const_element_type* input_A
,const_dim_type input_A_ndim, const_dim_type* input_A_dims
,const_element_type* input_B
,const_dim_type input_B_ndim, const_dim_type* input_B_dims
,element_type* output_Y
,const_dim_type output_Y_ndim, const_dim_type* output_Y_dims
```

```
8   + #include "${OperatorName}.hh"
9   + #include "diagnostic/msgHandling.hh"
10  + #include <Eigen/Dense>
11  +
12  + namespace cimHW {
13  +
14  + cimHW${OperatorName}Op::cimHW${OperatorName}Op(
        const std::string &configFile
        ${argument_list}
      )
        : cimHWOp("${OperatorName}")
19  + , m_cimCU(configFile)
20  + {
21  +   verbose1(opName());
22  + }
23  +
24  + void cimHW${OperatorName}Op::simulate()
25  + {
26  +   // copy back to output
27  +   //memcpy(m_output_Y, result.data(), sizeof(numType) * ;
28  + }
29  +
30  + } // namespace cimHW
```

ONNC

# Create a new operator class from the template

- op.template.cc

```
 8  + #include "${OperatorName}.hh"
 9  + #include "diagnostic/msgHandling.hh"
10  + #include <Eigen/Dense>
11  +
12  + namespace cimHW {
13  +
14  + cimHW${OperatorName}Op::cimHW${OperatorName}Op(
15  +    const std::string &configFile
16  +      ${argument_list}
17  +    )
18  +        : cimHWOp("${OperatorName}")
19  +    , m_cimCU(configFile)
20  + {
21  +    verbose1(opName());
22  + }
23  +
24  + void cimHW${OperatorName}Op::simulate()
25  + {
26  +    // copy back to output
27  +    //memcpy(m_output_Y, result.data(), sizeof(numType) * ;
28  + }
29  +
30  + } // namespace cimHW
```

```
// map to proper view of matrix according to transision
m_matrix_A = Eigen::Map<const MatrixXfRowMajor>(input_A, input_A_dims[0], input_A_dims[1]);
m_matrix_B = Eigen::Map<const MatrixXfRowMajor>(input_B, input_B_dims[0], input_B_dims[1]);
```

ONNC

# Create a new operator class from the template

- op.template.cc

```
23  +
24  + void cimHW${OperatorName}Op::simulate()
25  + {
26  +         // take matrixA as kernel and matrixB as data to use ComputeUnit
27  +    //me   MatrixXfRowMajor result = MatrixXfRowMajor::Zero(m_matrix_A.rows(), m_matrix_B.cols());
28  + }           for(int i=0; i<m_matrix_A.rows(); i++)
29  +                 result.row(i) = m_cimCU.compute(m_matrix_B, m_matrix_A.row(i).transpose());
30  + } // n
              // copy back to output
              memcpy(m_output_Y, result.data(), sizeof(numType) * m_matrix_A.rows() * m_matrix_B.cols());

              verbose1("\n");
```

ONNC

# Register to desired factory class

- skysim/onnc-cimHW/lib/factory/cimHWOpFactory.hh
  - Add declaration
  - replace ${argument_list} skysim/onnc-cimHW/include/abstractOpFactory.hh
    - 2nd parameter of create${OperatorName}Op function

```
49  +  // virtual std::unique_ptr<cimHWOp> create${OperatorName}Op(
50  +  //    void* context
51  +  //    ${argument_list}
52  +  // ) override;
```

```
,const_element_type* input_A
,const_dim_type input_A_ndim, const_dim_type* input_A_dims
,const_element_type* input_B
,const_dim_type input_B_ndim, const_dim_type* input_B_dims
,element_type* output_Y
,const_dim_type output_Y_ndim, const_dim_type* output_Y_dims
```

ONNC

# Register to desired factory class

- skysim/onnc-cimHW/lib/factory/cimHWOpFactory.cc
  - Add include

```
 7  + //===----------------------------------------------------------
 8  + #include "cimHWOpFactory.hh"
 9  + #include "../conv.hh"
10  + #include "../gemm.hh"
11  +                                  #include "../matmul.hh"
12  + namespace cimHW {
13  +
```

ONNC

# Register to desired factory class

- skysim/onnc-cimHW/lib/factory/cimHWOpFactory.cc
  - Add definition
    - replace ${OperatorName
    - replace ${argument_list}
    - remove type in 2nd ${argument_list}

```
77  +
78  + // std::unique_ptr<cimHWOp> cimHWOpFactory::create${OperatorName}Op(
79  + //    void* context
80  + //    ${argument_list}
81  + // )
82  + // {
83  + //    const std::string* cimCon:
84  + //    return std::make_unique<c:
85  + //       *cimConfig
86  + //       ${argument_list}
87  + //    );
88  + // }
89  +
90  + } // namespace cimHW
```

```
,const_element_type* input_A
 ,const_dim_type input_A_ndim, const_dim_type* input_A_dims
 ,const_element_type* input_B
 ,const_dim_type input_B_ndim, const_dim_type* input_B_dims
 ,element_type* output_Y
 ,const_dim_type output_Y_ndim, const_dim_type* output_Y_dims
```

```
, input_A, input_A_ndim, input_A_dims
, input_B, input_B_ndim, input_B_dims
, output_Y,output_Y_ndim, output_Y_dims
```

ONNC

# Adding the new file to building system

- Add the matmul.cc to skysim/onnc-cimHW/lib/CMakeLists.txt

# Rebuild

- $ cd /onnc/onnc-cimHW.build
  $ make;
  $ sudo make install


- Run a test case to verify the implementation

- $ ./example/matmul-layer share/cimConfig.json

ONNC

Open Neural Network Compiler