# TMVA SOFIE
## Enhancing the Machine Learning Inference Engine

Sanjiban Sengupta
EP-SFT

**Supervisor:** Lorenzo Moneta

/sanjibansg        /sanjiban-sengupta

- Toolkit for Multivariate Analysis
- Provides a Machine Learning environment for training, testing and evaluation of multivariate methods.

SOFIE

# SOFIE

System for

# SOFIE

System for Optimized

# SOFIE

System for Optimized Fast

# SOFIE

## System for Optimized Fast Inference

# SOFIE

**S**ystem for **O**ptimized **F**ast Inference code **E**mit

# SOFIE

System for Optimized Fast Inference code Emit

*inference code, fast to operate, with least dependencies*

Motivation

## Motivation

- ML ecosystem mostly focuses on model training.

## Motivation

- ML ecosystem mostly focuses on model training.
- Inference in Tensorflow & PyTorch
  - supports only their own model
  - usage of C++ environment is difficult
  - heavy dependency

## Motivation

- ML ecosystem mostly focuses on model training.
- Inference in Tensorflow & PyTorch
  - supports only their own model
  - usage of C++ environment is difficult
  - heavy dependency
- Inference in ONNX (Open Neural Network Exchange)
  - can use ONNXRuntime by Microsoft
  - large dependency
  - difficult to integrate in HEP applications
    - control of libraries, threads
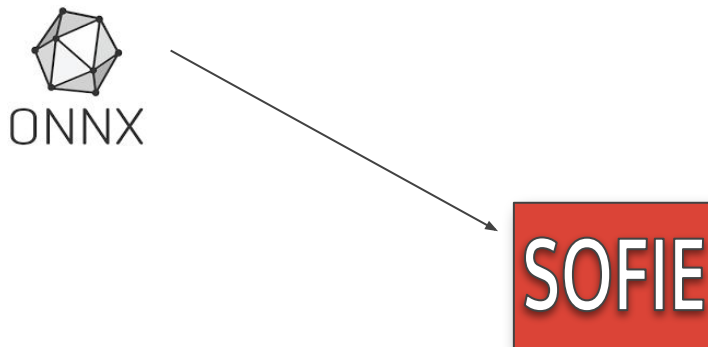    - not optimized for single event evaluation

- Intermediate representation following ONNX standards.

- Intermediate representation following ONNX standards.
- Inference code generation with least latency and minimal dependency

- Intermediate representation following ONNX standards.
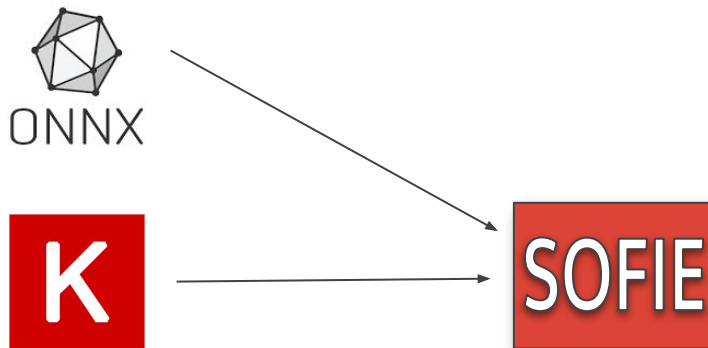- Inference code generation with least latency and minimal dependency

- Intermediate representation following ONNX standards.
- Inference code generation with least latency and minimal dependency
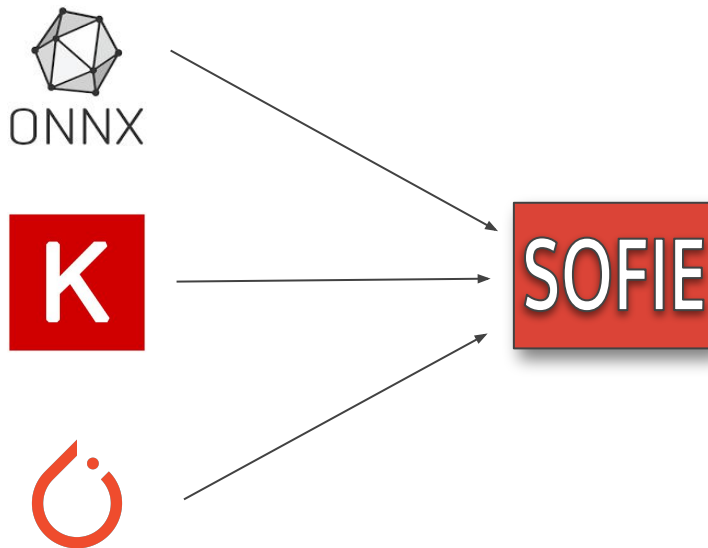
- Intermediate representation following ONNX standards.
- Inference code generation with least latency and minimal dependency

- Intermediate representation following ONNX standards.
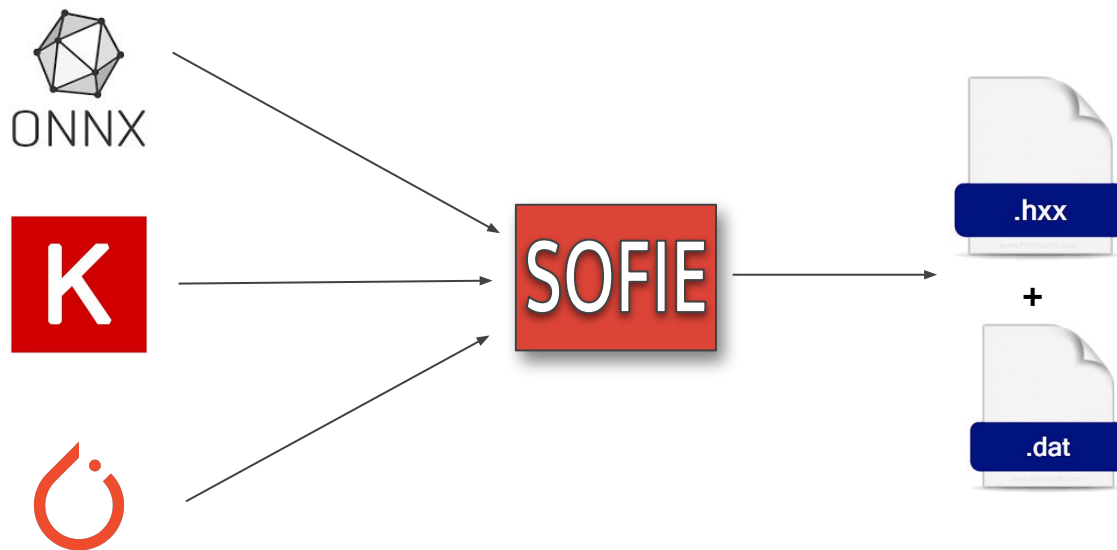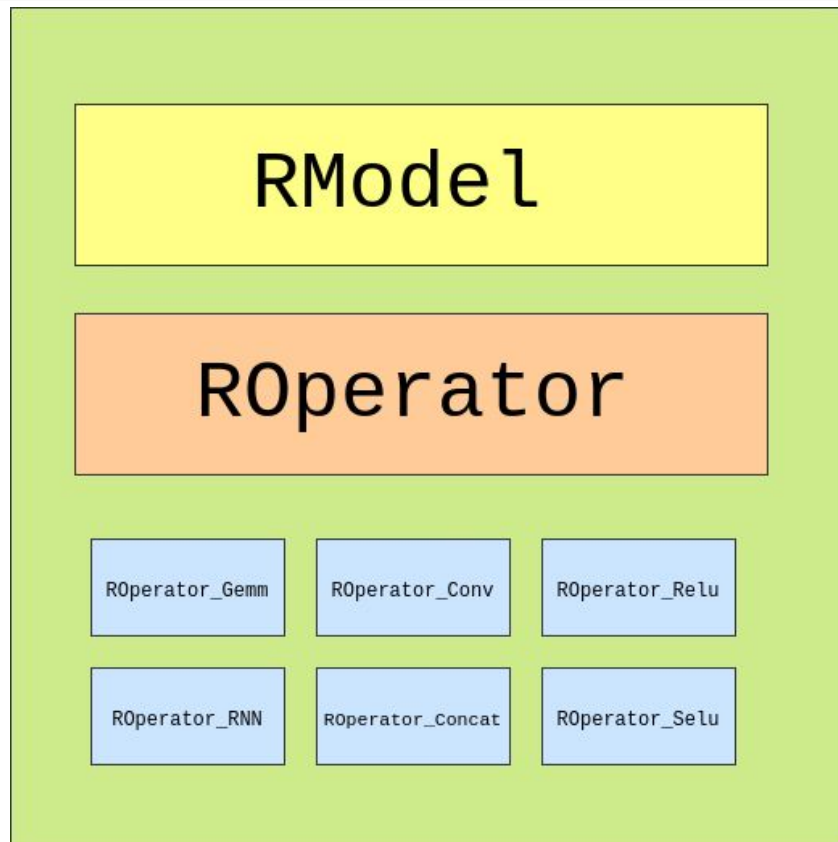- Inference code generation with least latency and minimal dependency

Parser for translating an ONNX model to SOFIE's IR

```cpp
using namespace TMVA::Experimental::SOFIE;
RModelParser_ONNX parser;
RModel model = parser.Parse("model.onnx");
```

Parser for translating an ONNX model to SOFIE's IR

```
using namespace TMVA::Experimental::SOFIE;
RModelParser_ONNX parser;
RModel model = parser.Parse("model.onnx");
```

Parser for translating Keras (.h5) models

```
SOFIE::RModel model = SOFIE::PyKeras::Parse("KerasModel.h5");
```

Parser for translating an ONNX model to SOFIE's IR

```cpp
using namespace TMVA::Experimental::SOFIE;
RModelParser_ONNX parser;
RModel model = parser.Parse("model.onnx");
```

Parser for translating Keras (.h5) models

```cpp
SOFIE::RModel model = SOFIE::PyKeras::Parse("KerasModel.h5");
```

Inference code generation

```cpp
// generate text code internally (with some options)
model.Generate();
// write output header file and data weight file
model.OutputGenerated();
```

## SOFIE's Generated code



```cpp
// Code auto generated by TMVA SOFIE

namespace TMVA_SOFIE_Linear_event{

struct Session {

Session(std::string filename ="") {
 if (filename.empty()) filename = "Linear_event.dat";
 std::ifstream f;
 f.open(filename);
 // read weight data file
 ...................
}
std::vector<float> infer(float* tensor_input1){
 ...................
```

- Extending support of SOFIE Keras parser

- Extending support of SOFIE Keras parser
- Implement SOFIE Custom operator support

- Extending support of SOFIE Keras parser
- Implement SOFIE Custom operator support
- Implement support for parsing Graph Neural Networks in SOFIE

- No native support for ONNX translation

- No native support for ONNX translation
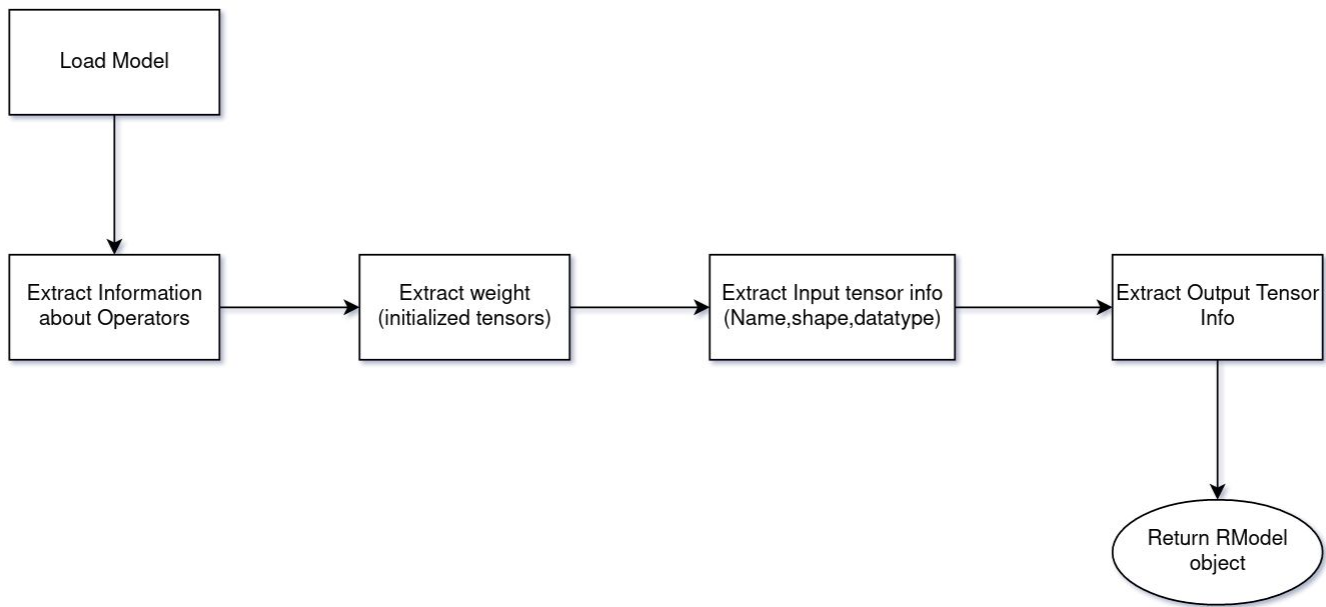- TF2ONNX may convert a Keras .h5 model to ONNX

- No native support for ONNX translation
- TF2ONNX may convert a Keras .h5 model to ONNX
- SOFIE Keras Parser!
  - simpler to use
  - no need for input spec
  - built on latest opset

## Algorithm for Parser

```
┌──────────────┐
│  Load Model  │
└──────┬───────┘
       │
       ▼
┌──────────────────┐    ┌──────────────────┐    ┌──────────────────────┐    ┌──────────────────────┐
│ Extract          │───▶│ Extract weight   │───▶│ Extract Input tensor │───▶│ Extract Output Tensor│
│ Information      │    │ (initialized     │    │ info                 │    │ Info                 │
│ about Operators  │    │ tensors)         │    │ (Name,shape,datatype)│    │                      │
└──────────────────┘    └──────────────────┘    └──────────────────────┘    └──────────┬───────────┘
                                                                                        │
                                                                                        ▼
                                                                              ╭───────────────────╮
                                                                              │   Return RModel   │
                                                                              │      object       │
                                                                              ╰───────────────────╯
```

## Current Support

| Keras Layer | Status |
| --- | --- |
| Dense | Implemented & Integrated |
| Permute | Implemented & Integrated |
| ReLU, Selu, Sigmoid | Implemented & Integrated |
| Batch Normalization | PR Merged |
| Convolution (2D) | PR Merged |
| Reshape | PR Merged |
| Basic Binary Operators: Add, Subtract, Multiply | PR Under Review |
| Activations: Softmax, LeakyRelu, Tanh | PR Drafted |
| Concat | PR Drafted |

- ONNX standards specifies 183 operators currently.

- ONNX standards specifies 183 operators currently.
- Need a custom user operator specification

- ONNX standards specifies 183 operators currently.
- Need a custom user operator specification
  - simple to define

# SOFIE Custom Operator

- ONNX standards specifies 183 operators currently.
- Need a custom user operator specification
  - simple to define
  - easy to test, debug, and evaluate

# SOFIE Custom Operator

- ONNX standards specifies 183 operators currently.
- Need a custom user operator specification
  - simple to define
  - easy to test, debug, and evaluate
  - few overheads and dependencies
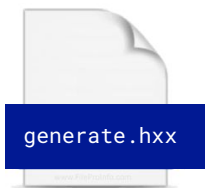
Definition

## Definition

- Define custom operator with the required attributes
  - Operator name
  - Input tensor names
  - Output tensor names
  - Output tensor shapes
  - Header file name

## Definition

- Define custom operator with the required attributes
  - Operator name
  - Input tensor names
  - Output tensor names
  - Output tensor shapes
  - Header file name
- Define Compute function in Header file

| generate.hxx | weights.dat | compute.hxx |

## Definition

- Define custom operator with the required attributes
    - Operator name
    - Input tensor names
    - Output tensor names
    - Output tensor shapes
    - Header file name
- Define Compute function in Header file

- High demands of Graph Neural Networks in High energy physics research

- High demands of Graph Neural Networks in High energy physics research
- CMS

- High demands of Graph Neural Networks in High energy physics research
- CMS
  - uses Particlenet; graph neural network supporting graph convolution, i.e. edge convolution and dynamic graph CNN methods

- High demands of Graph Neural Networks in High energy physics research
- CMS
  - uses Particlenet; graph neural network supporting graph convolution, i.e. edge convolution and dynamic graph CNN methods
- LHCb

- High demands of Graph Neural Networks in High energy physics research
- CMS
  - uses Particlenet; graph neural network supporting graph convolution, i.e. edge convolution and dynamic graph CNN methods
- LHCb
  - plans to use DeepMind's Graph Nets library; builds GNN on top of Tensorflow & Sonnet

- High demands of Graph Neural Networks in High energy physics research
- CMS
  - uses Particlenet; graph neural network supporting graph convolution, i.e. edge convolution and dynamic graph CNN methods
- LHCb
  - plans to use DeepMind's Graph Nets library; builds GNN on top of Tensorflow & Sonnet

Current Plans & Implementation

## Current Plans & Implementation

- Following the DeepMind's Graph Nets architecture

## Current Plans & Implementation
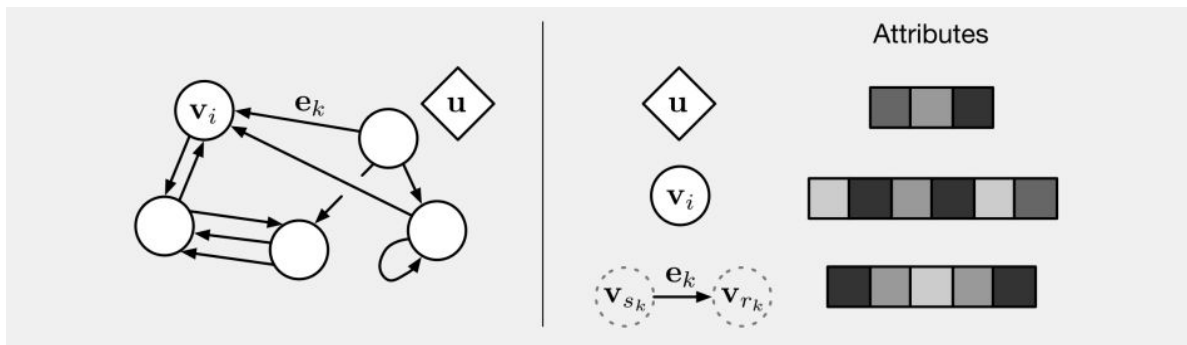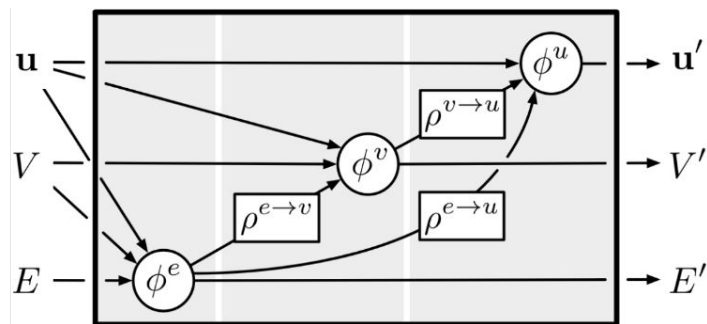
- Following the DeepMind's Graph Nets architecture
- Intermediate Representation
    - Nodes list
    - Edges list
    - Sender's list
    - Receiver's list
    - Global values

## Current Plans & Implementation

- Following the DeepMind's Graph Nets architecture
- Intermediate Representation
  - Nodes list
  - Edges list
  - Sender's list
  - Receiver's list
  - Global values

## Current Plans & Implementation

- Following the DeepMind's Graph Nets architecture
- Intermediate Representation
  - Nodes list
  - Edges list
  - Sender's list
  - Receiver's list
  - Global values

## Current plans for implementation

- Following the DeepMind's Graph Nets architecture
- Intermediate Representation
  - Nodes list
  - Edges list
  - Sender's list
  - Receiver's list
  - Global values

- Operating functions
  - Updation functions
  - Aggregation functions

- Finishing implementation & integration of Graph Neural Network support in SOFIE

- Finishing implementation & integration of Graph Neural Network support in SOFIE
- Implementing support for new operators in TMVA SOFIE

# Conclusion

- Link to Forked Repository
  github.com/sanjibansg/root

- Link to SOFIE in current ROOT master
  github.com/root-project/root/tree/master/tmva/sofie

- Link to TMVA/SOFIE tutorials
  root.cern.ch/doc/master/group__tutorial__tmva.html

- Link to SOFIE notebooks
  github.com/lmoneta/tmva-tutorial/tree/master/sofie

## Using SOFIE's Generated code

```
// Code generated automatically by TMVA for Inference of Model file [model.h5] at [Wed Aug  3 20:32:37 2022]

#include "Model.hxx"
// create session class
TMVA_SOFIE_Model::Session s();
//–- event loop
…….
{
// evaluate model: input is an array of type float *
std::vector<float> result = s.infer(input);
}
```

## Using SOFIE's Generated code

```python
import ROOT
# compile generate SOFIE code using ROOT interpreter
ROOT.gInterpreter.Declare('#include "Model.hxx"')
# create session class
s = ROOT.TMVA_SOFIE_Model.Session()
#-- event loop
…….
 # evaluate the model , input can be a numpy array of type float32
 result = s.infer(input)
```
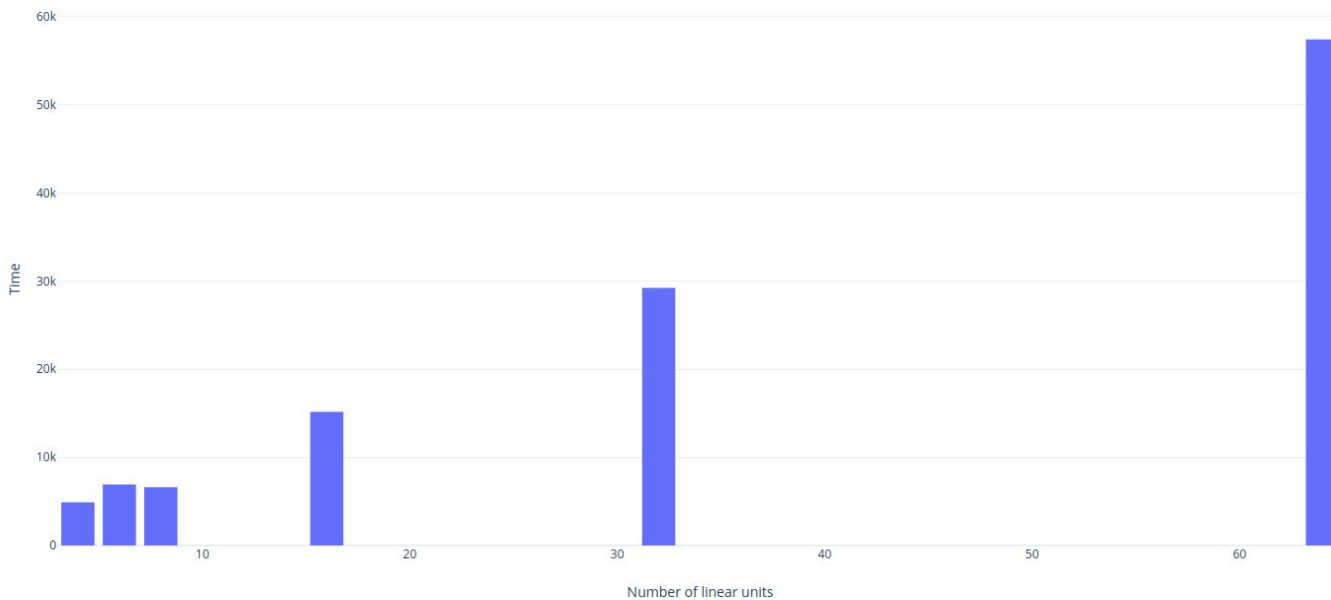
## Definition

- Define custom operator with the required attributes
  - Operator name
  - Input tensor names
  - Output tensor names
  - Output tensor shapes
  - Header file name
- Define Compute function in Header file

## Interface

```
std::unique_ptr<SOFIE::ROperator> op;
op.reset(new SOFIE::ROperator_Custom<float>("Exp", {"denseBiasAdd0"}, {"exp_out"}, {{1,4}},
"exp_compute.hxx"));
```
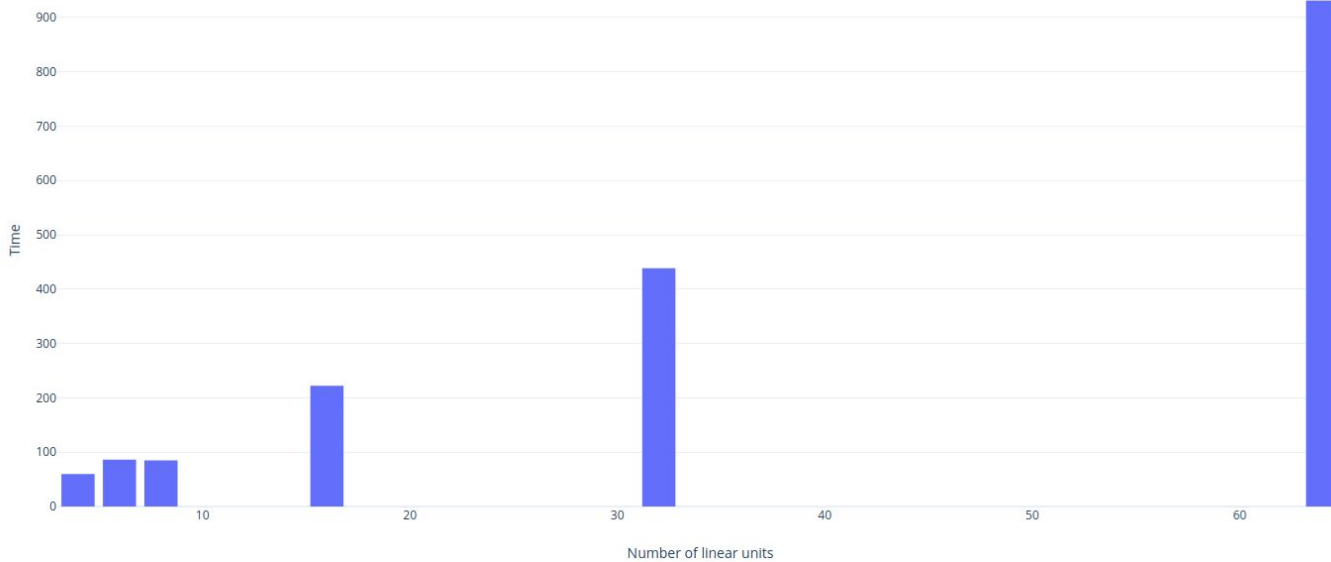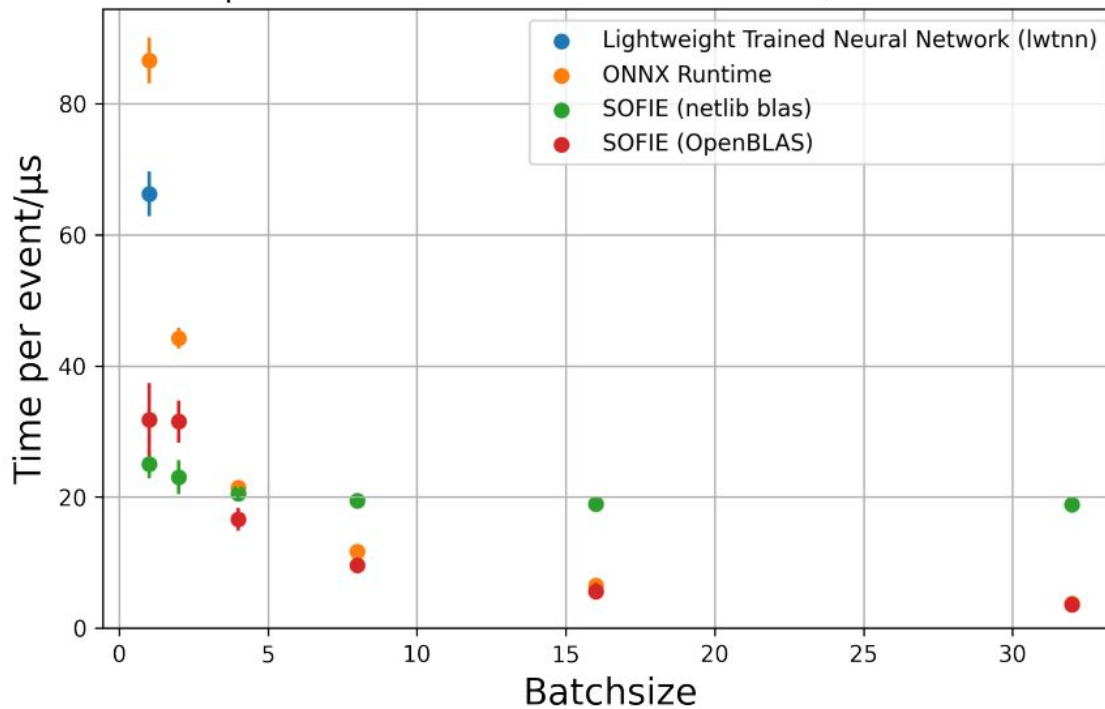
## Benchmarking for Keras

## Benchmarking for SOFIE
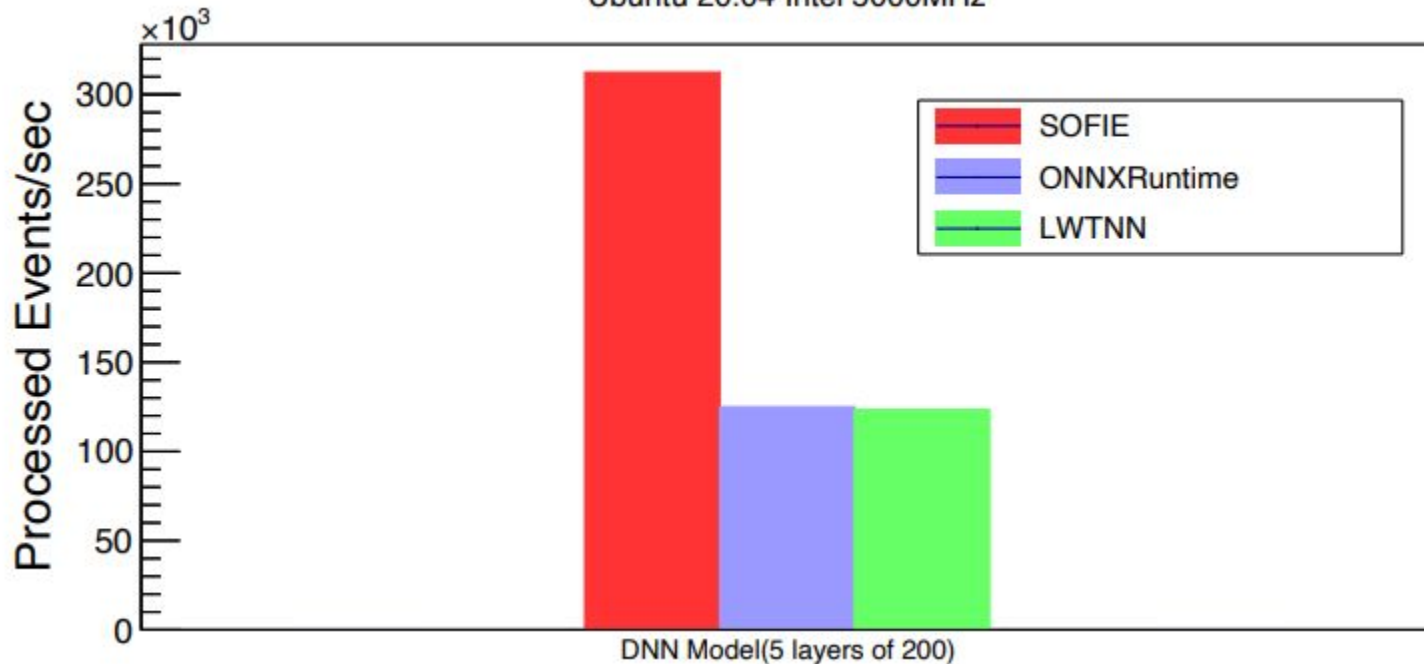
Time per event for different batch size, cache flushed

Ubuntu 20.04 Intel 5000MHz

**Larger = Better**

Processed Events/sec

SOFIE
ONNXRuntime
LWTNN

DNN Model(5 layers of 200)

Ubuntu 20.04 Intel 5000MHz (Batch Size = 1)