

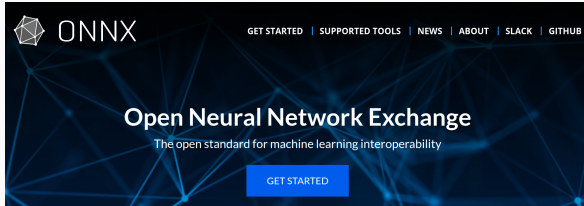
# ONNX and its usage in Geant4

---

Thomas Madlener

Feb 04, 2022

# Open Neural Network Exchange - ONNX



 [onnx/onnx](https://github.com/onnx/onnx)  
[onnx.ai](https://onnx.ai)

- Open ecosystem providing an open source format for AI models
  - Originated from Facebook (as Toffee), now also supported by e.g. Microsoft and Amazon (and many others)
- Provides the definition of an extensible computation graph model, built-in operators and standard data types
  - A common intermediate representation (IR) of the computation graphs employed by the different ML frameworks
  - IR spec is versioned to allow for evolution
- Allows to transfer models between different ML frameworks
  - Choose the one that is best suited for the task at hand

- Open source engine with the ability to run ONNX models
  - Training and inference
  - Supported on multiple platforms for different hardware
- Bindings for several languages available
  - Python, C++, C#, Java, JavaScript, Julia,...
- Deeply integrated with the whole Microsoft (Azure) stack<sup>1</sup>
- Shared library with C++ API available in LCG software stack
  - And used for inference in `extended/parameterisations/Par04` of `geant4`



---

<sup>1</sup>at least if you believe their website ;)

# geant4 Par04 example

- Cylindrical calorimeter setup with the possibility to use an ML model for fast simulation energy deposits
  - Support for ORT and [LWTNN](#)
  - See [source](#) for some more details
- Comes with pre-trained model
  - VAE with latent space  $\in \mathbb{R}^{10}$ , with i.i.d.  $\sim \text{No}(0, 1)$  and 4 condition variables (energy, angle + one-hot encoding of two different geometries)
  - Energy and angle labels are normalized to training range
  - Outputs a 1D vector of energy deposits on a 3D mesh (mesh coincides with mesh used for full simulation training samples)
- Rotates 3D mesh into place depending on incident particle momentum
- [src/Par04InferenceSetup.cc](#) for details

# Running the example - environment setup

- Use latest (and now only) HEAD release: `HEAD-2022-01-30`
- Needs some additional environment setup to pull in ORT from `LCG_101` release
  - I couldn't get it to run with the setup that is used in the README!
  - Use [setup\\_env.sh](#) from `onnx-integration` branch of [MarlinML](#) that does a complete setup (including iLCSoft HEAD release)

```
git clone -b onnx-integration https://gitlab.desy.de/ilcsoft/MarlinML
source MarlinML/setup_env.sh
```

# Running the example - compiling & running

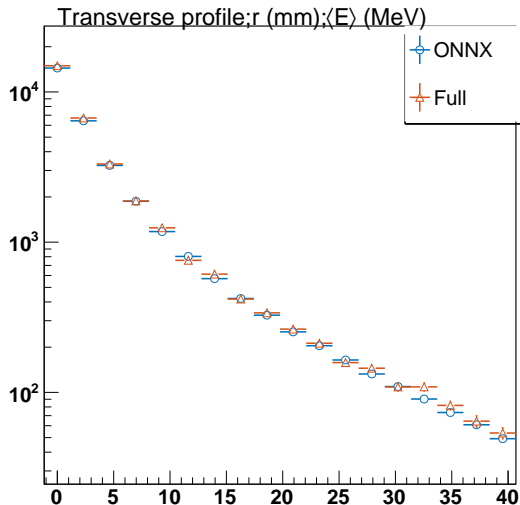
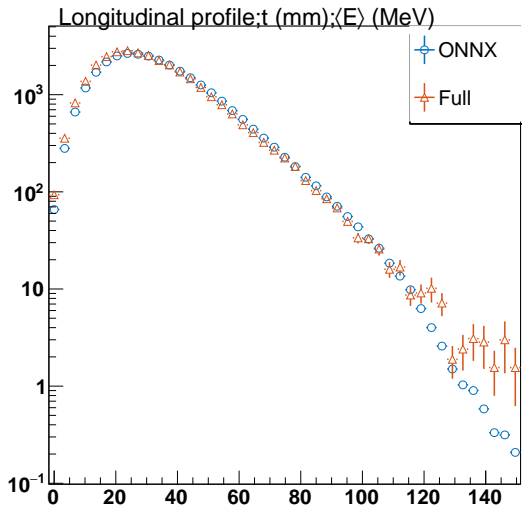
- Copy example sources and compile

```
cp -r ${Geant4_DIR}/../../examples/extended/parameterisations/Par04 .
cd Par04
mkdir build && cd build
cmake ..
make -j4
```

- There should be a message `ONNX Runtime inference library found.`
  - Potentially also a warning that can be ignored
- Run (according to README) full simulation and ONNX

```
./examplePar04 -m examplePar04.in # full simulation
./examplePar04 -m examplePar04_onnx.in # fast simulation w/ ONNX
```

# Compare Full and FastSim



- Default settings shoot 10 GeV electron along y-axis

# Plotting script

```
import ROOT as r
from utils.plot_helpers import mkplot
from utils.setup_plot_style import set_basic_style
set_basic_style()
r.gStyle.SetStatStyle(0) # remove stats box


r.gROOT.SetBatch() # no display, just plots

f_onnx = r.TFile.Open('10GeV_100events_fastsim_onnx.root')
f_full = r.TFile.Open('10GeV_100events_fullsim.root')

keys = [k.GetName() for k in f_onnx.GetListOfKeys()]
hnames = filter(lambda k: f_onnx.Get(k).InheritsFrom('TH1D'), keys)

for name in hnames:
    can = mkplot([f_onnx.Get(name), f_full.Get(name)],
                 legEntries=['ONNX', 'Full'], legPos=(0.78, 0.78, 0.98, 0.95))
    can.SaveAs(f'{name}.pdf')

    can = mkplot([f_onnx.Get(name), f_full.Get(name)],
                 legEntries=['ONNX', 'Full'], legPos=(0.78, 0.78, 0.98, 0.95),
                 logy=True)
    can.SaveAs(f'{name}_log.pdf')
```

- `utils` are from  [tmadlener/chib\\_chic\\_polFW](https://github.com/tmadlener/chib_chic_polFW)  
Get via `svn co https://github.com/tmadlener/chib_chic_polFW/trunk/python/utils utils`



# ORT in MarlinML

- Added an example processor using ORT to MarlinML
  - [onnx-integration branch](#)
- Example of how to export PyTorch model to ONNX and use it in inference in a Marlin processor

```
17 class ExampleModule(torch.nn.Module):
18     """Very simple example module that doesn't do anything useful"""
19
20     def __init__(self):
21         super(ExampleModule, self).__init__()
22         self.linear = torch.nn.Linear(4, 4)
23
24     def forward(self, x, h):
25         vals = torch.tanh(self.linear(x) + h)
26         return vals
27
28     @staticmethod
29     def rand_inputs():
30         """Get random inputs for the module for cases where the values are
31            but the shape information is important (e.g. ONNX export)"""
32         return (torch.rand(3, 4), torch.rand(3, 4))
33
34
35     def _store_model_torch(model, filename):
36         """Store the model as torch script file"""
37         script_model = torch.jit.script(model)
38         script_model.save(filename)
39
40
41     def _store_model_onnx(model, filename):
42         """Store the model in the onnx format"""
43         inputs = model.rand_inputs()
44         # Very basic here. The export method has many more params!
45         torch.onnx.export(model, inputs, filename,
46                           # For testing some C++ API details
47                           input_names=('x', 'h'),
48                           output_names=('o',)
49                           )
```

```

147 // ( 344C_L_1 = 0, 1 ~ 11111111, 771 ) 1
148 // Setup the input tensors
149 //
150 // NOTE: most of these should be const, but the API doesn't like that too
151 // much for some reason
152 /*const*/ auto x =
153     m_linearInput ? linspace(0.f, 1.f, 12) : randomVals(0.f, 1.f, 12);
154 /*const*/ auto h =
155     m_linearInput ? linspace(1.f, 0.f, 12) : randomVals(0.f, 1.f, 12);
156
157 const std::vector<int64_t> inputDims = {dim1, dim2};
158 // Here we have to go in a slightly roundabout way as well by first
159 // creating the tensors and then pushing them into the vector as it seems
160 // to not be possible to directly initialize the vector for some reason...
161 /*const*/ auto xTensor = Ort::Value::CreateTensor<float>(
162     m_info, x.data(), dim1 * dim2, inputDims.data(), inputDims.size());
163 /*const*/ auto hTensor = Ort::Value::CreateTensor<float>(
164     m_info, h.data(), dim1 * dim2, inputDims.data(), inputDims.size());
165
166 std::vector<Ort::Value> inputs;
167 inputs.push_back(std::move(xTensor));
168 inputs.push_back(std::move(hTensor));
169
170 auto outputs = m_session->Run(Ort::RunOptions{nullptr}, m_inputNames.data(),
171     inputs.data(), inputs.size(),
172     m_outputNames.data(), m_outputNames.size());
173
174 // Here we know that we only have one output variable from the model, so
175 // can directly use that
176 const auto &result = outputs.front();
177 // printTypeInfo("result:", result.GetTypeInfo(), "outputs[0]", 42);
178
179 const auto values =
180     std::vector<float>(result.GetTensorData<float>()),
181     result.GetTensorData<float>() + dim1 * dim2);

```

## Inference Examples/examples

## Marlin\_run\_example\_onnx\_model.xml

```

MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.80460, 0.72127, 0.91316, 0.391818, 0.168772, 0.553464, 0.602551, 0.8023113, 0.474475, 0.508492, 0.850201, 0.026433)
Run 0:
x = ( 0.94442, 0.02339, 0.403504, 0.437769, 0.725374, 0.383536, 0.430949, 0.0312443, 0.481316, 0.399449, 0.184888)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "k = ( 0.705337, 0.031816, 0.574109, 0.323207, 0.897250, 0.067070, 0.31446, 0.723601, 0.280812, 0.642674, 0.561815)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "h = ( 0.702384, 0.313681, 0.400119, 0.877805, 0.905517, 0.440765, 0.249112, 0.121011, 0.314843, 0.605781, 0.52844, 0.080933)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 0:
x = ( 0.961309, 0.348838, 0.802933, 0.382331, 0.323259, 0.017333, 0.155349, 0.335344, 0.118467, 0.08723, 0.073120, 0.559137)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.320146, 0.023744, 0.852739, 0.139239, 0.272469, 0.559448, 0.760387, 0.0403013, 0.708407, 0.577036, 0.46155, 0.653093)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "k = ( 0.530555, 0.105366, 0.405449, 0.630981, 0.515922, 0.452434, 0.210182, 0.009111, 0.211218, 0.462111, 0.021104, 0.5555)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 0:
x = ( 0.72444, 0.959712, 0.422384, 0.51438, 0.132159, 0.724115, 0.724115, 0.400071, 0.50645, 0.45645, 0.46932, 0.925464)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "h = ( 0.32028, 0.49021, 0.0247992, 0.5423, 0.01845, 0.468688, 0.0189151, 0.360505, 0.032923, 0.313660, 0.50805, 0.212762)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.80774, 0.72104, 0.24732, 0.434019, 0.20393, 0.34464, 0.17479, 0.50629, 0.527118, 0.546442, 0.42007, 0.32797)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 0:
x = ( 0.440824, 0.355091, 0.304677, 0.327546, 0.240573, 0.38813, 0.791469, 0.205160, 0.181001, 0.317531, 0.207606, 0.950270)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.47444, 0.74444, 0.74444, 0.74444, 0.74444, 0.74444, 0.74444, 0.74444, 0.74444, 0.74444, 0.74444, 0.74444)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 1:
x = ( 0.738941, 0.385158, 0.203643, 0.446304, 0.367287, 0.0129083, 0.431318, 0.916473, 0.359987, 0.68712, 0.489991, 0.523953)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "h = ( 0.8074793, 0.107185, 0.448169, 0.201629, 0.122923, 0.188084, 0.427712, 0.83813, 0.512044, 0.5895, 0.442392, 0.241292)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.82318, 0.493504, 0.322384, 0.349297, 0.468232, 0.254201, 0.371866, 0.765151, 0.394466, 0.32502, 0.187465, 0.313023)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 2:
x = ( 0.607112, 0.154793, 0.807541, 0.35552, 0.4678, 0.30536, 0.230463, 0.40222, 0.491939, 0.164743, 0.402304, 0.728693)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "h = ( 0.791509, 0.154435, 0.766113, 0.08309554, 0.246434, 0.404451, 0.51146, 0.200363, 0.226862, 0.077944, 0.90588, 0.116668)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.937554, 0.831645, 0.430979, 0.345057, 0.027762, 0.37197, 0.22344, 0.20847, 0.777271, 0.55443, 0.76833, 0.216044)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 0:
x = ( 0.804607, 0.0446191, 0.405556, 0.143411, 0.904341, 0.80022573, 0.6708883, 0.12044, 0.434468, 0.193937, 0.808209, 0.0661440)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "h = ( 0.548526, 0.46167, 0.348783, 0.794659, 0.092811, 0.892348, 0.260232, 0.545713, 0.355136, 0.92323, 0.4825, 0.728164)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.803155, 0.296434, 0.813209, 0.839249, 0.819381, 0.872204, 0.094874, 0.719283, 0.871748, 0.406614, 0.156231, 0.687147)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 0:
x = ( 0.451104, 0.506795, 0.209226, 0.970742, 0.503472, 0.029076, 0.164446, 0.245493, 0.299235, 0.4796, 0.710484, 0.350870)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.50872, 0.953891, 0.508011, 0.607351, 0.276011, 0.37746, 0.07746, 0.0078088, 0.287405, 0.209124, 0.209124, 0.209123)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "k = ( 0.869083, 0.002149, 0.402113, 0.648399, 0.039161, 0.31318, 0.270339, 0.193467, 0.740915, 0.219362, 0.400772, 0.178765)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 5:
x = ( 0.49328, 0.05463, 0.89441, 0.714015, 0.971919, 0.718083, 0.314783, 0.104482, 0.104482, 0.59139, 0.20168, 0.460602)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "h = ( 0.64852, 0.175309, 0.404023, 0.140601, 0.262991, 0.341482, 0.451793, 0.129935, 0.371914, 0.596393, 0.0200874, 0.4494933)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.82318, 0.493504, 0.322384, 0.349297, 0.468232, 0.254201, 0.371866, 0.765151, 0.394466, 0.32502, 0.187465, 0.313023)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 0:
x = ( 0.2317, 0.913248, 0.109947, 0.090919, 0.0031846, 0.100812, 0.100812, 0.100812, 0.100812, 0.100812, 0.100812, 0.100812)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "h = ( 0.154047, 0.280291, 0.093115, 0.564134, 0.0104, 0.63129, 0.407580, 0.940213, 0.181349, 0.770779, 0.166569, 0.564911)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.70375, 0.377096, 0.609831, 0.588085, 0.441326, 0.61243, 0.330210, 0.73241, 0.87046, 0.37044, 0.202424, 0.374091)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 0:
x = ( 0.24113, 0.738232, 0.15375, 0.815881, 0.84132, 0.416673, 0.718603, 0.287128, 0.244667, 0.469775, 0.814469, 0.199935)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "h = ( 0.30455, 0.40859, 0.7044, 0.42180, 0.974478, 0.909349, 0.909349, 0.909349, 0.909349, 0.909349, 0.909349, 0.909349)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.71914, 0.801749, 0.09118, 0.121449, 0.792949, 0.792949, 0.792949, 0.792949, 0.792949, 0.792949, 0.792949, 0.792949)
MESSAGE "MyOnnxRuntimeExampleProcessor"
Run 0:
x = ( 0.852305, 0.250005, 0.40450, 0.341719, 0.581343, 0.870037, 0.319374, 0.770055, 0.809192, 0.522426, 0.834754, 0.175060)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "h = ( 0.7352, 0.472248, 0.704549, 0.911893, 0.653897, 0.137742, 0.386804, 0.896726, 0.10281, 0.714907, 0.80394, 0.369393)
MESSAGE "MyOnnxRuntimeExampleProcessor"
MESSAGE "y = ( 0.6147, 0.107188, 0.511797, 0.728466, 0.924466, 0.150795, 0.8242449, 0.747474, 0.872097, 0.814645, 0.718109, 0.483091)

```