# ONNX based inference in ATLAS

ILC software meeting
*https://agenda.linearcollider.org/event/10344/*

15th May 2024

Antonio Giannini

# *About me*

**A. Giannini, antonio.giannini@cern.ch**

- PhD degree at Universita' Federico II, Naples (Italy)

- Post-doc at USTC, Hefei (China)

- Area of interest

    - Diboson physics (BSM searches and VBS measurements)

    - Boosted objects tagging, W/Z and Hbb/cc

    - ML applications, event classification, background estimations, jet tagging, Anomaly detection

- ATLAS responsibilities

    - HDBS Physics group ML liaison 2021 - 2023

    - Conveners of Diboson searches sub-group (HDBS-DBL), 2022 - 2024

*HDBS = Higgs and DiBoson Searches*
*DBL = DiBoson Lab*

# Outline

**ONNX based inference in ATLAS**

- Introduction

  ▷ typical workflow

  ▷ why inference?

- ML inference in ATLAS

  ▷ historical evolution, from lwtnn to ONNX

- Applications overview

  ▷ (lwtnn) for analysis level

  ▷ ONNX for jet tagging

*ATHENA software (public)*
*documentation*

**Wrap-up**

*Let me remember to cite the proper*
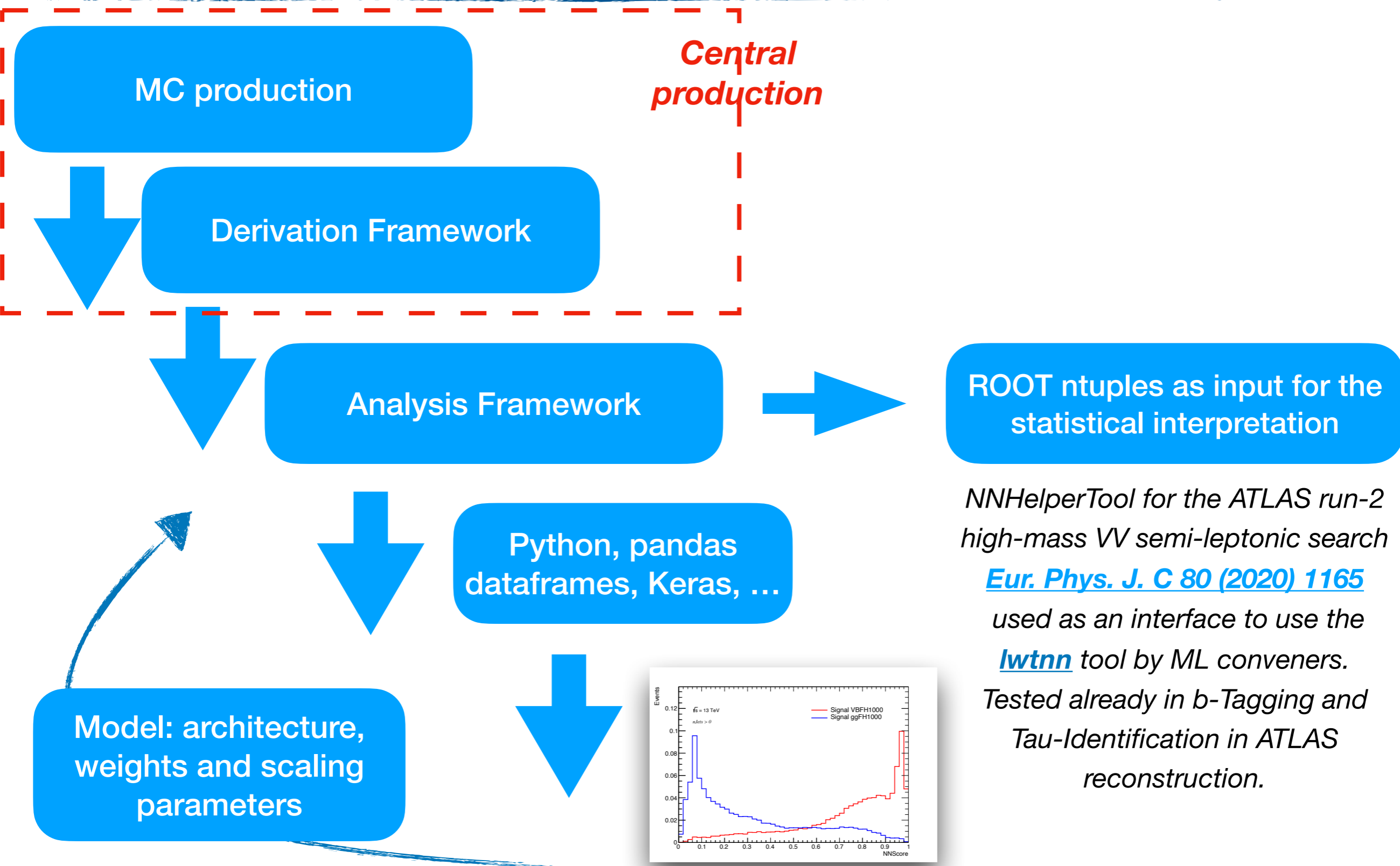*reference for discussing the ATLAS software!*
*ATL-SOFT-PUB-2021-001.pdf*

# ML models vs analyses

- *ML applications have had a huge boost in Physics analyses at collider experiments*
  - during my ATLAS experience, ~O(1) in 2017 up to ~O(10) in 2024 of ML applications in analyses
  - searches better use cases for applications/development than measurements
- *Wide variety of python based frameworks are used now*
  - TensorFlow, Keras, PyTorch, sci-kit learn, etc etc
  - relatively easy to get started, many tutorials around, many people willing to pass the knowledge
- This implies
  - a large population of "custom" frameworks doing mostly the same things
  - a large population of datasets, intermediate/output formats, ROOT, pandas, h5, etc
- *All of these ingredients can make really hard the reproducibility of results!!!*

# Typical workflow: entry level

MC production

*Central production*

Derivation Framework

Analysis Framework

ROOT ntuples as input for the statistical interpretation

Python, pandas dataframes, Keras, …

Model: architecture, weights and scaling parameters



*NNHelperTool for the ATLAS run-2 high-mass VV semi-leptonic search **Eur. Phys. J. C 80 (2020) 1165** used as an interface to use the **lwtnn** tool by ML conveners. Tested already in b-Tagging and Tau-Identification in ATLAS reconstruction.*

# Typical workflow: next level

MC production

*Central production*

Derivation Framework, with dedicated information

"Analysis" Framework

Python, pandas dataframes, Keras, …

Model: architecture, weights and scaling parameters
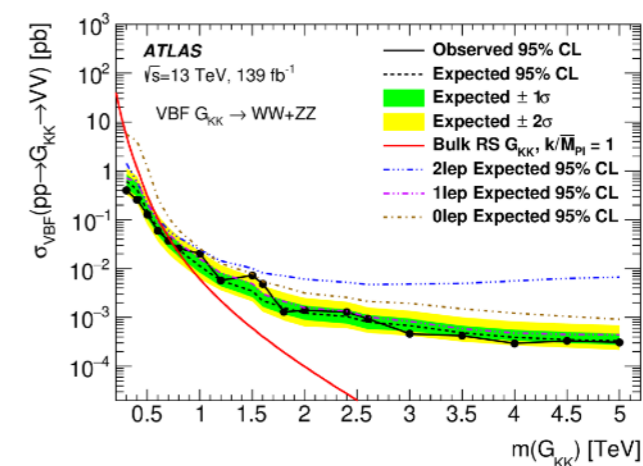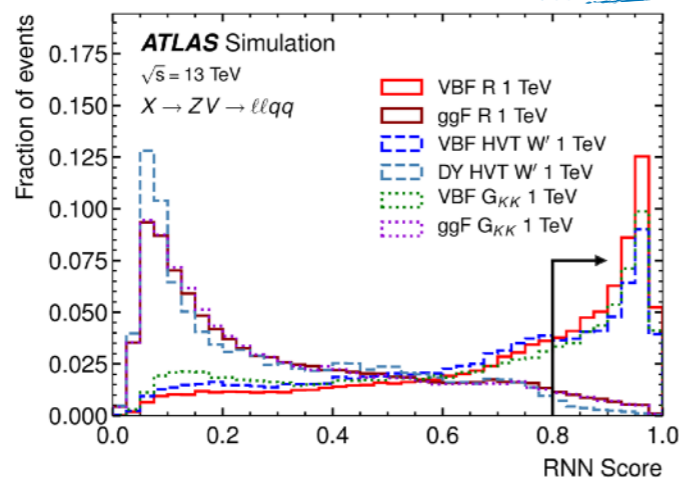
*Central production*

MC production

Derivation Framework with light information

Analysis Framework

ROOT ntuples as input for the statistical interpretation

# *Why ML models inference in C++*

- ML models inference in C++ is crucial for two main aspects

- ***Reproducibility of results outside the collaboration***

    ▷ event classifier (NN>x), final discriminant (to complete the analysis

    ingredients)

    ▷ anomaly detections (when the ML is the main analysis ingredient…)

- ***Deployment of a tool that can be used for several analyses***

    ▷ jet tagging, low-level variables based classifiers, etc

- ***Inference is not a plus, nor an aesthetic element, it is a significant part of the***

    ***R&D of a ML tool for physics!***

# ML inference in ATLAS

- Two main tools have been used so far
  - ▷ *lwtnn: many analyses applications in the latest few years*
  - ▷ *ONNX: appeared in the last 3-4 years, so far, mainly jet tagging approaches (b-tag, Hbb/cc, W/Z) have been implemented*
- An historical perspective
  - ▷ my memories about this start around 2016/17 or so, at the time we had lwtnn tool available and it saved my PhD thesis!
  - ▷ found a funny definition from Dan Guest (former ATLAS ML convener, lwtnn author) *"Dark ages"* before lwtnn was made



**Back in the dark ages... (Sep 2015)**

- There were no inference engines
- Every ML framework
  - Had its own serialization
  - Did its own inference        *Dan talk 2022*
  - Had a million dependencies
  - Had zero stability
- So we wrote our own (lwtnn)
  - Serialized to JSON
  - Separate **inputs, outputs, data, graph**
  - Started small: sequential models
    - Latter added graphs
- We store networks on a file system
  - Write only, locally cached on demand

random Tau ID model I found

5

- In 2018/19 when I implemented the inference for my first analysis only few (lwtnn) use cases were available:
  - ▷ tau RNN-ID, b-tagging, 1 SUSY multi jets analysis

# lwtnn implementation

- lwtnn has been widely used in ATLAS
  - quite user friendly, you can get familiar with it in an afternoon and wrap it in your favourite tool/analysis code
- It is available in ATLAS software but you can also install it on your environment
- Model conversion support, very human readable!

**Supported Layers**

*lwtnn github*

In particular, the following layers are supported as implemented in the Keras sequential and functional models:

| | K sequential | K functional |
|---|---|---|
| Dense | yes | yes |
| Normalization | See Note 1 | See Note 1 |
| Maxout | yes | yes |
| Highway | yes | yes |
| LSTM | yes | yes |
| GRU | yes | yes |
| Embedding | sorta | issue |
| Concatenate | no | yes |
| TimeDistributed | no | yes |
| Sum | no | yes |

*paper HEPData*
*VBF-RNN architecture .json*

```
1  {  "class_name": "Model",
2     "config": { "name": "model_1",
3                    "layers": [{"name": "input_1",
4                                 "class_name": "InputLayer",
5                                 "config": { "batch_input_shape": [null, 2, 4],
6                                             "dtype": "float32",
7                                             "sparse": false,
8                                             "name": "input_1"
9                                 },
10                   "inbound_nodes": []
11                   },
12                   {"name": "jet_masking",
13                    "class_name": "Masking",
14                    "config": { "name": "jet_masking",
15                                "trainable": true,
16                                "mask_value": -99.0},
17                    "inbound_nodes": [[["input_1", 0, 0, {}]]]},
18                    {  "name": "jet_lstm1",
19                       "class_name": "LSTM",
20                       "config": { "name": "jet_lstm1",
21                                   "trainable": true,
22                                   "batch_input_shape": [null, null, 4],
23                                   "dtype": "float32",
24                                   "return_sequences": true,
25                                   "return_state": false,
26                                   "go_backwards": false,
27                                   "stateful": false,
28                                   "unroll": false,
29                                   "units": 25,
30                                   "activation": "tanh",
31                                   "recurrent_activation": "hard_sigmoid",
32                                   "use_bias": true,
```

# ONNX in a nutshell

- ONNX runtime (**[https://github.com/onnx](https://github.com/onnx)**)
  - ▷ more an industry based usage
  - ▷ based on **Open Neural Network eXchange** files
  - ▷ supports many more operators than lwtnn (factor 10 maybe)
  - ▷ but not universal support
  - ▷ it is a bit less user friendly than lwtnn
- Installation made in the ATLAS software (ATHENA), ONNX is available as an external dependency
  - ▷ **External/onnxruntime**
  - ▷ for ATLAS user: no need to install anything, just use it!!!

- lwtnn

  - tau RNN-ID (***ATL-PHYS-PUB-2019-033***)

  - b-tagging RNN/DeepSets (***ATL-PHYS-PUB-2020-014***)

  - VBF-RNN tagger in VV semi-leptonic (***Eur. Phys. J. C 80 (2020) 1165***)

  - Anomaly detection search (YXh) in X(qq)h(bb) events (no routine but available in HEPData) (***Phys. Rev. D 108 (2023) 052009***)

- ONNX

  - Multi b-jets + MET SUSY search (***Eur. Phys. J. C 83 (2023) 561***)

  - b-tagging (***ATL-PHYS-PUB-2022-027***)

  - Hbb/cc boosted tagger (***ATL-PHYS-PUB-2023-021***)

# *Analysis level application example (lwtnn)*

- Run model conversion (keras —> .json)

- Features scaling config file supported

- Simple tool to interface the physics input
  preparation and load the ML model



```json
{} neural_net.json    80.54 KiB
 1  {
 2    "defaults": {},
 3    "inputs": [
 4      {
 5        "name": "NJets",
 6        "offset": -3.393768,
 7        "scale": 0.660834
 8      },
 9      {
10        "name": "Zcand_m",
11        "offset": -90.953651,
12        "scale": 0.210485
13      },
14      {
15        "name": "Zcand_pt",
16        "offset": -626.260071,
17        "scale": 0.00172
18      },
19      {
20
21
22
```

**Features scaling**

**(Physics) inputs preparation**

**Retrieve the ML model**

**CxAODFramework/ Run2PaperVersion//NNHelper.cxx**

```cpp
601  lwt::VectorMap NNHelper::get_JetVars(const std::vector<lwt::Input>& input
602    lwt::VectorMap out, out_mask;
603    // ramp through the input multiplier
604    const size_t total_inputs = inputs.size(); // nVariables = [pT, eta, ph
605    for (size_t nnn = 0; nnn < total_inputs; nnn++) {
606      const auto& input = inputs.at(nnn);
607      out[input.name] = {};
608      out_mask[input.name] = {};
609
610      double mean = 0., std = -1.;
611
612      if( input.name=="variable_0" ){
613        if(UseMax6Jets) out.at(input.name) = { Jet1_pt , Jet2_pt , Jet3_pt
614        if(UseMax2Jets) out.at(input.name) = { Jet1_pt , Jet2_pt };
615        mean = mean_pT;
616        std  = std_pT;
617      }
```
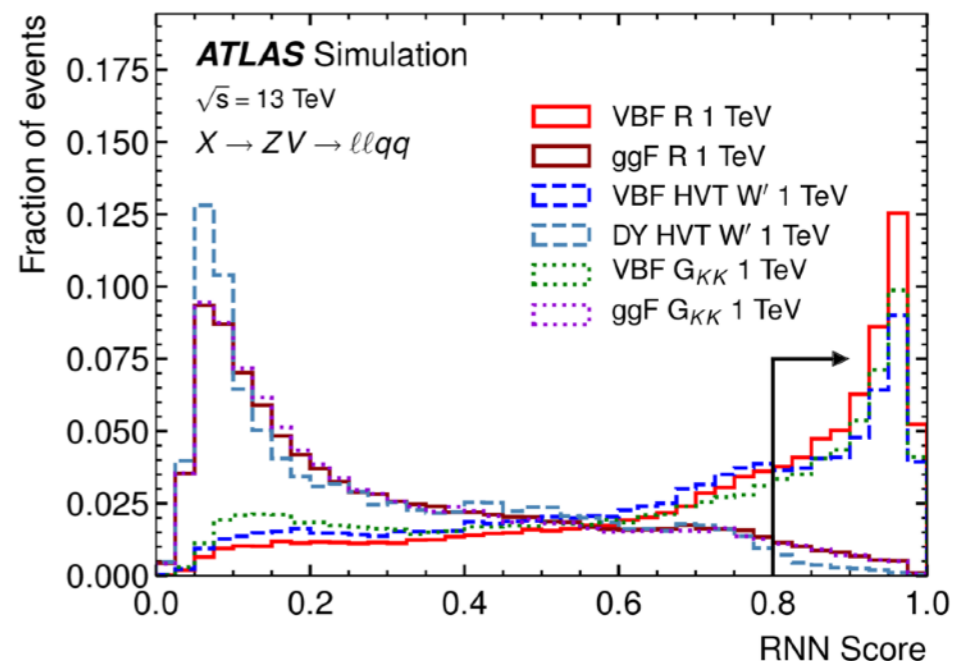
```cpp
 91  void NNHelper::LoadRNN(){
 92
 93    if(IsDebug) std::cout << " --- LoadRNN ---" << std::endl;
 94
 95    lwt::LightweightGraph tagger(config, config.outputs.begin()->first);
 96
 97    LightweightGraph::SeqNodeMap seq = get_sequences(config.input_sequences);
 98    LightweightGraph::NodeMap in_nodes;
 99    //in_nodes = get_generated(config.inputs);
100
101    // Loop over the output names and compute the output for each
102    for (const auto& output: config.outputs) {
103      auto out_vals = tagger.compute(in_nodes, seq, output.first);
104      if(IsDebug){
105        std::cout << output.first << ":" << std::endl;
106        for (const auto& out: out_vals) {
107          std::cout << out.first << " " << out.second << std::endl;
108        }
109      }
110      RNNScore = out_vals["out_0"];
111      if(IsDebug) std::cout << RNNScore << "   " << Jet1_pt << "    " << Jet1_e
112    }
```
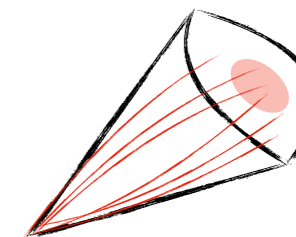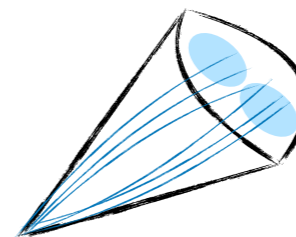
# *Jet tagger application: W/Z tagger*

- Historical Boosted Jet Tagger (BJT) tool with implementation of many W/Z/top jet taggers
    - ▷ *atlas/athena/ BoostedJetTaggers*
    - ▷ cut based taggers, lwtnn DNN, ONNX CNN/DNN
    - ▷ constituents based W/Z tagger has been implemented using ONNX
- Same config structure as for baseline taggers,
    - ▷ just more info are added as needed

```
35
36   DecorationName: SmoothWContained50DisCoJet
37
38   pTCutLow: 200.0
39
40   pTCutHigh: 3000.0
41
42   MassCutLow: (7.180e+01)+(1.124e-02)*pow(x,1)+(-1.721e-05)*pow(x,2)+(4.113e-09)*pow(x,3)
43
44   MassCutHigh: (1.255e+02)+(-7.448e-02)*pow(x,1)+(5.402e-05)*pow(x,2)+(-1.118e-08)*pow(x,3)
45
46   DisCoDNNCut: (6.731e-01)+(3.434e-05)*pow(x,1)+(1.595e-08)*pow(x,2)+(-7.777e-12)*pow(x,3)
47
48   CNNTaggerFileName: CNNTaggers/ONNXModel_VVJJ_Model_m175RScaling_18feb22/model.onnx
49   nPixelsEta: 40
50   aEta: -1.
51   bEta: 1.
52   nPixelsPhi: 40
53   aPhi: -1.
54   bPhi: 1.
55   nColors: 3
56   DoRScaling: 1
57   ### the jet pT is in MeV
58   RScaling_p0: 0.0761
59   RScaling_p1: 379.e+3
60
61   ### DisCo model
62   DisCoTaggerFileName: CNNTaggers/ONNXModel_VVJJ_Model_DisCoJet_13mar22/model.onnx
63
64   ### features scaling
65   pT_mean: 1002977.7292303614
66   pT_std: 626014.5212713333
67   CNN_mean: 0.44037336334870836
68   CNN_std: 0.25781137380904184
69   D2_mean: 1.8977217190151305
70   D2_std: 1.1563601945608166
71   nTracks_mean: 32.33996946809005
72   nTracks_std: 12.06663413423996
```

**Usual Config file**

**pT-parametrisations**

**CNN model**

**DisCo-DNN model**

# First constituents based tagger: ONNX helper

- Dedicated class to interface the "physics" data to the ONNX usage
  - **_JSSMLTool.cxx_**
  - inspired from the main ATHENA example, **_AthenaExamples/ AthExOnnxRuntime_**
- ONNX usage:
  - use the shared ATHENA service **_AthOnnxruntimeService_**
  - inputs in Ort::Value format
- So far, implemented DNN based and CNN based architectures



*(Physics) inputs preparation*

*Retrieve the ML model*

```
195    //preparing container to hold input data
196
197    size_t input_tensor_size = m_nbins_eta*m_nbins_phi*m_ncolors;
198    std::vector<float> input_tensor_values(input_tensor_size);
199
200    int testSample = 0;
201    input_tensor_values = input_tensor_values_[testSample];
202
203    //preparing container to hold output data
204    int output_tensor_values = output_tensor_values_[testSample];
205
206    // create input tensor object from data values
207    auto memory_info = Ort::MemoryInfo::CreateCpu(OrtArenaAllocator, OrtMemTypeDefault);
208    Ort::Value input_tensor = Ort::Value::CreateTensor<float>(memory_info, input_tensor_va
209    assert(input_tensor.IsTensor());
210
211    auto output_tensors = m_session->Run(Ort::RunOptions{nullptr}, m_input_node_names.data
212    assert(output_tensors.size() == 1 && output_tensors.front().IsTensor());
213
```

# Jet tagger application example: Hbb/cc tagger

- Implementation in the FlavorTagging ATLAS groups

- Supported taggers (based on GNN/transformers!)

  ▷ *small-R jets b-/c-tagging*

  ▷ *large-R jets Xbb/cc tagging*

- More advanced architecture already available

  ▷ more work done in the metadata and handling
    of the tool

  ▷ you can get nice inspiration!

*BTagging/BTagConfig.py*



*FlavorTagDiscriminants/GNN.cxx*

```
69     ],
       "AntiKt10UFOCSSKSoftDropBeta100Zcut10": [
70         "BTagging/20230413/gn2xv00/antikt10ufo/network.onnx",
71         "BTagging/20230413/gn2xwithmassv00/antikt10ufo/network.onnx",
72         "BTagging/20230705/gn2xv01/antikt10ufo/network.onnx",
73     ],
```

*model version
configuration*

```
236     # Associate tracks to the jet
237     result.merge(JetParticleAssociationAlgCfg(
238         inputFlags,
239         jetcol,
240         trackCollection,
241         JetTrackAssociator,
242     ))
```

*(Physics) inputs preparation*

*Extract metadata
from the ONNX model*

```
22    GNN::GNN(const std::string& nn_file, const GNNOptions& o):
23      m_onnxUtil(nullptr),
24      m_jetLink(jetLinkName),
25      m_defaultValue(o.default_output_value)
26    {
27      // Load and initialize the neural network model from the given file path.
28      std::string fullPathToOnnxFile = PathResolverFindCalibFile(nn_file);
29      m_onnxUtil = std::make_shared<OnnxUtil>(fullPathToOnnxFile);
30
31      // Extract metadata from the ONNX file, primarily about the model's inputs.
32      auto lwt_config = m_onnxUtil->getLwtConfig();
33
34      // Create configuration objects for data preprocessing.
35      auto [inputs, constituents_configs, options] = dataprep::createGetterConfig(
36          lwt_config, o.flip_config, o.variable_remapping, o.track_link_type);
37
```

# *First GNN based tagger: ONNX helper*

**FlavorTagDiscriminants/GNN.cxx**

- Implementation in the FlavorTagging ATLAS groups
- Also in this case the core part is using ONNX session
    - ▷ after physics data preparation into tensor data
    - ▷ inputs prepared in Ort::Value format

***Retrieve the ML model***

```
150    std::vector<float> input_tensor_values;
151
152    // create input tensor object from data values
153    auto memory_info = Ort::MemoryInfo::CreateCpu(
154      OrtArenaAllocator, OrtMemTypeDefault
155    );
156    std::vector<Ort::Value> input_tensors;
157    for (auto const &node_name : m_input_node_names){
158      input_tensors.push_back(Ort::Value::CreateTensor<float>(
159        memory_info, gnn_inputs.at(node_name).first.data(), gnn_inputs.at(node_name).fi
160        gnn_inputs.at(node_name).second.data(), gnn_inputs.at(node_name).second.size())
161      );
162    }
163
164    // casting vector<string> to vector<const char*>. this is what ORT expects
165    std::vector<const char*> input_node_names;
166    input_node_names.reserve(m_input_node_names.size());
167    for (const auto& name : m_input_node_names) {
168      input_node_names.push_back(name.c_str());
169    }
170    std::vector<const char*> output_node_names;
171    output_node_names.reserve(m_output_nodes.size());
172    for (const auto& node : m_output_nodes) {
173      output_node_names.push_back(node.name_in_model.c_str());
174    }
175
176    // score model & input tensor, get back output tensor
177    // Although Session::Run is non-const, the onnx authors say
178    // it is safe to call from multiple threads:
179    //   https://github.com/microsoft/onnxruntime/discussions/10107
180    Ort::Session& session ATLAS_THREAD_SAFE = *m_session;
181    auto output_tensors = session.Run(Ort::RunOptions{nullptr},
182      input_node_names.data(), input_tensors.data(), input_node_names.size(),
183      output_node_names.data(), output_node_names.size()
184    );
185
```

# *The talk in one slide*

| | lwtnn (jet tagging, flavour tagging, analyses) | ONNX (jet tagging) | ONNX (flavour tagging) |
|---|---|---|---|
| **available (examples)** | yes | | |
| **ATHENA links** | *BoostedJetTaggers/ JSSWTopTaggerDNN.cxx* | *AthOnnxRuntimeBJT/ JSSMLTool.cxx* | *FlavorTagDiscriminants/ GNN.cxx* |
| **model converter** | supported | need an easy script | |
| **supported operators** | ~20 | many more (but not universal) | |
| **model metadata** | human friendly .json | need to decode the information | |
| **data pre-paration** | we are the physicists! | | |
| **feed the inputs** | config for features scaling | similar but need to take care of the scaling | |
| **tool configuration** | depends by the case | more user level configuration (easy to implement your custom tagger!) | more central configuration |
| **architectures exploited so far** | DNN/RNN/etc | DNN/CNN | DNN/GNN |
| **ATLAS fw step** | Derivation/Analysis level | Analysis level | Derivation level |

# Conclusions

**ONNX based inference in ATLAS**

- Inference is a crucial point in ML tool developments
  - ▷ production of a specific tool (1 jet tagger —> many analyses )
  - ▷ reproducibility of physics results
- Since the "dark age" people have extensively started to use tools for inference
  - ▷ lwtnn mainly, migrating to ONNX now
  - ▷ ONNX supported as ATLAS external, wider operators support
- ONNX example code application in the (public) ATHENA software

*Thanks for listening, hope it can be useful for you!!!*

*One ring to rule all of them*



*One inference tool to run all the analyses*



*Do it! Without inference it will be really hard to reproduce your fancy ML model!*

*backup*

# ATHENA software

https://atlassoftwaredocs.web.cern.ch/athena/athena-intro/

## Athena Introduction

Last update: 06 Mar 2024 [History] [Edit]

Athena is based on the common Gaudi framework that is used by ATLAS, LHCb and FCC.

Athena code is hosted in the CERN GitLab service where the repository owned by the `atlas` user plays a special role as the repository from which all production releases will be built. In the ATLAS development workflow this repository is also the merge target for developments from individual users or groups. This repository includes a number of different *branches* (defined as "an independent line of development") being actively committed to at any given time, which may be synchronized to varying degrees (either manually or via automated procedures) depending on their purpose. There is much more information about the development workflow in the ATLAS git development tutorial.

The main projects are:

| Project | Purpose |
|---|---|
| Athena | Reconstruction and Derivation production[*] |
| AthGeneration | For event generation |
| AthSimulation | For full Geant4 simulation |
| AthAnalysisBase | Athena based analysis |
| AnalysisBase | Non-athena ROOT based analysis |
| DetCommon | For reading trigger configuration when e.g. configuring L1 Hardware |

# Motivation: Production ML in HEP

## HEP Land

- Lots of C++
- Event Pipelines
- Lots of legacy
- Moves slow

## ML Land

- Python
- Batched
- Minimal legacy
- Moves fast

- Many Dependencies

2

# What do need to retrieve? And how? ONNX runtime

|  | **CNN** | **(DisCo)DNN** |
|---|---|---|
| **Data pre-processing** | ImageBuilding | Access to JSS variables |
| **Inputs** | 3x40x40 pixels | N variables after features scaling |
| **model** | via ONNX | via ONNX |

**CNNScore**

**DisCoJet is made by 2 scores**

*Low-level, jet constituents (taste-aware network)*

**DNNScore**
**(Mass de-correlated via DisCo)**

**Jet D2**

*High-Level Jet variables*

**Jet nTracks**

**Tagger WP(s)**

**Jet Mass**

**Other O(15) JSS variables**

***DisCo reference and implementation here: arxiv.2001.05310**

# DisCoJet in the BJT: implementation (I)

**Initialise the tool:**
- we need help (*AthExOnnxRuntime_BJT*) to interface to ONNX session
- exploit the nice configuration file feature to set whatever parameter you would need

**[Retrieving NN scores]**

```
110    if(m_LoadCNNTagger){
111        // init tool
112        std::string CNNTaggerFileName = m_configReader.GetValue("CNNTaggerFileName", "aaa");
113        std::string m_ModelPath = ("/data/BoostedJetTaggers/SmoothedWZTaggers/"+CNNTaggerFileName);
114
115        ATH_MSG_INFO("SmoothedWZDisCoJetTagger::MLBosonTagger()" << "   + ModelPath " << m_ModelPath );
116        m_MLBosonTagger = new AthONNX::CxxApiAlgorithm("MLBosonTagger", NULL, m_ModelPath);
117        ATH_CHECK( m_MLBosonTagger -> initialize() );
118
119        // get model paramters from the config file
120        ATH_MSG_INFO("SmoothedWZDisCoJetTagger::MLBosonTagger() read value from config" );
121
122        // set parameters
123        ATH_CHECK( m_MLBosonTagger->setProperty("nPixelsEta", (int)m_configReader.GetValue("nPixelsEta", -99)) );
124        ATH_CHECK( m_MLBosonTagger->setProperty("aEta", m_configReader.GetValue("aEta", -99.)) );
125        ATH_CHECK( m_MLBosonTagger->setProperty("bEta", m_configReader.GetValue("bEta", -99.)) );
126        ATH_CHECK( m_MLBosonTagger->setProperty("nPixelsPhi", (int)m_configReader.GetValue("nPixelsPhi", -99)) );
127        ATH_CHECK( m_MLBosonTagger->setProperty("aPhi", m_configReader.GetValue("aPhi", -99.)) );
128        ATH_CHECK( m_MLBosonTagger->setProperty("bPhi", m_configReader.GetValue("bPhi", -99.)) );
129        ATH_CHECK( m_MLBosonTagger->setProperty("nColors", m_configReader.GetValue("nColors", -99)) );
130        ATH_CHECK( m_MLBosonTagger->setProperty("DoRScaling", (bool)m_configReader.GetValue("DoRScaling", -99)) );
131        ATH_CHECK( m_MLBosonTagger->setProperty("RScaling_p0", m_configReader.GetValue("RScaling_p0", -99.)) );
132        ATH_CHECK( m_MLBosonTagger->setProperty("RScaling_p1", m_configReader.GetValue("RScaling_p1", -99.)) );
133
134        m_MLBosonTagger -> PrintInfo();
135
136    }
```

```
490    // use ML tool on constituents
491    m_MLBosonTagger -> MakeJetImage("Charged" , jet, csts_charged );
492    m_MLBosonTagger -> MakeJetImage("Neutral" , jet, csts_neutral );
493    m_MLBosonTagger -> MakeJetImage("Combined", jet, csts_combined);
494
495    // evaluate the model
496    ATH_CHECK( m_MLBosonTagger -> execute() );
497
498    // save ML score
499    CNNScore = m_MLBosonTagger -> GetScore();
500    //   std::cout << " --- CNNScore, " << CNNScore << std::endl;
501
502    // reset images and go to next jet
503    m_MLBosonTagger -> ResetImages();
504
505    return CNNScore;
```

*Actions for the CNN:*
- build images
- load the score

*Good practise!*
*check ALWAYS that the C++ implementation is matching (within float precision) your result from keras!*

**Loop jet by jet:**

```
269    // Retrieve the CNN score
270    float jet_cnn (-99.);
271    float jet_disco (-99.);
272
273    if(m_LoadCNNTagger){
274        jet_cnn = GetCNNScore(&jet);
275        (*m_dec_cnn)(jet) = jet_cnn;
276    }
277    if(m_UseDisCoTagger){
278        auto scaler = ReadScaler(); // this shoud be moved in the init
279        auto JSSVars = GetJSSVars(jet, scaler);
280        jet_disco = GetDisCoDNNScore(jet, JSSVars);
281        (*m_dec_disco)(jet) = jet_disco;
282    }
```

**Retrieve the info you need:**
- CNN Score
- DNN Score

```
509    float SmoothedWZDisCoJetTagger::GetDisCoDNNScore(const xAOD::Jet& jet, co
510
511        // init value
512        float DisCoDNNScore (-99.);
513
514        // set inputs
515        //auto JSSVars = GetJSSVars(jet);
516        //std::map<std::string, double> JSSVars = GetJSSVars(jet);
517
518        m_MLBosonTagger_DisCo -> SetJSSInputs(JSSVars);
519
520        // evaluate the model
521        ATH_CHECK( m_MLBosonTagger_DisCo -> executeDisCo() );
522
523        // save ML score
524        DisCoDNNScore = m_MLBosonTagger_DisCo -> GetDNNScore();
525        //   std::cout << " --- CNNScore, " << CNNScore << std::endl;
526
527        return DisCoDNNScore;
528
529    }
```

*Actions for the DNN:*
- retrieve all JSS vars
- load the score

# DisCoJet in the BJT: implementation (II)

**[Actual tagger WP]**

```
175    ATH_MSG_INFO( "Smoothed WZ Tagger tool initialized" );
176    ATH_MSG_INFO( "  Mass cut low      : " << m_strMassCutLow );
177    ATH_MSG_INFO( "  Mass cut High     : " << m_strMassCutHigh );
178    if( m_UseCNNTagger )
179      ATH_MSG_INFO( "  CNN cut low       : " << m_strCNNCut );
180    else if( m_UseDisCoTagger )
181      ATH_MSG_INFO( "  DisCoDNN cut low      : " << m_strDisCoCut );
182    else
183      ATH_MSG_INFO( "  D2 cut low        : " << m_strD2Cut );
184    if ( m_useNtrk )
185      ATH_MSG_INFO( "  Ntrk cut low      : " << m_strNtrkCut );
186    ATH_MSG_INFO( "  Decorate          : " << m_decorate );
187    ATH_MSG_INFO( "  DecorationName    : " << m_decorationName );
```

## Configuring the tagger:

- once all the variables are on the table, you just need to define the WP (i.e. cuts)
- example.1: UFO 3-Var Tagger
  - ▷ Mass, D2, nTracks
- example.2: UFO 3-Var CNN low-level Tagger
  - ▷ Mass, CNN, nTracks
- example.1: UFO 2-Var DisCoJet Tagger
  - ▷ Mass, DisCoJet

## Evaluate the cuts:

- the parametric cuts (as the pT) coming from the config file

```
284    /// Evaluate the values of the upper and lower mass bounds and the d2 cut
285    float cut_mass_low  = m_funcMassCutLow ->Eval(jet_pt);
286    float cut_mass_high = m_funcMassCutHigh->Eval(jet_pt);
287    float cut_d2 (-99.), cut_cnn (-99.), cut_disco (-99.);
288    if(m_UseCNNTagger)          cut_cnn      = m_funcCNNCut     ->Eval(jet_pt);
289    else if(m_UseDisCoTagger) cut_disco    = m_funcDisCoCut   ->Eval(jet_pt);
290    else                        cut_d2       = m_funcD2Cut      ->Eval(jet_pt);
```

```
302    /// Evaluate the cut criteria on mass and d2
303    if(m_UseCNNTagger){
304      ATH_MSG_VERBOSE( "Cut Values : MassWindow = [" << cut_mass_low << "," << cut_mass_
305      ATH_MSG_VERBOSE( "Cut Values : JetMass = " << jet_mass << ", CNN = " << jet_cnn )
306      if ( jet_cnn > cut_cnn ) m_accept.setCutResult( "PassCNN", true );
307    }
308    else if(m_UseDisCoTagger){
309      ATH_MSG_VERBOSE( "Cut Values : MassWindow = [" << cut_mass_low << "," << cut_mass_
310      ATH_MSG_VERBOSE( "Cut Values : JetMass = " << jet_mass << ", CNN = " << jet_disco
311      if ( jet_disco > cut_disco ) m_accept.setCutResult( "PassDisCoDNN", true );
312    }
313    else{
314      ATH_MSG_VERBOSE( "Cut Values : MassWindow = [" << cut_mass_low << "," << cut_mass_
315      ATH_MSG_VERBOSE( "Cut Values : JetMass = " << jet_mass << ", D2 = " << jet_d2 );
316      if ( jet_d2 < cut_d2 ) m_accept.setCutResult( "PassD2", true );
317    }
```

## Apply the cuts (tagger WP):

- apply the cut and store the results for the user!

**The full Tagger infrastructure has been reused!!!**

# DisCoJet in the BJT: running time

- The full infrastructure of the tool has been extended for this tagger
  - ▷ besides the specific variable or score, the actual 'Tagger' steps are always the same
- The user will not see any difference rather than in the specific print-outs and in the output results!

```
SmoothedContainedWTagg...INFO    Using config file : SmoothedContainedWTagger_AntiKt10VanillaSD_FixedSignalEfficiency50_2VarDisCoJet_v1.dat
SmoothedContainedWTagg...INFO    SmoothedWZTagger::MLBosonTagger()   + ModelPath /data/BoostedJetTaggers/SmoothedWZTaggers/CNNTaggers/ONNXModel_VVJJ_Model_m175RScaling_18feb22/model.onnx
PathResolver            WARNING Locating dev file dev/MLTest/2020-03-31/t10k-images-idx3-ubyte. Do not let this propagate to a release
PathResolver            WARNING Locating dev file dev/MLTest/2020-03-31/t10k-labels-idx1-ubyte. Do not let this propagate to a release
MLBosonTagger           INFO    Using model file: /srv/workDir/usr/WorkDir/0.0.1/InstallArea/x86_64-centos7-gcc8-opt//data/BoostedJetTaggers/SmoothedWZTaggers/CNNTaggers/
ONNXModel_VVJJ_Model_m175RScaling_18feb22/model.onnx
MLBosonTagger           INFO    Using pixel file: /cvmfs/atlas.cern.ch/repo/sw/database/GroupData/dev/MLTest/2020-03-31/t10k-images-idx3-ubyte
MLBosonTagger           INFO    Using pixel file: /cvmfs/atlas.cern.ch/repo/sw/database/GroupData/dev/MLTest/2020-03-31/t10k-labels-idx1-ubyte
MLBosonTagger           INFO    Created the ONNX Runtime session
MLBosonTagger           INFO    Output 0 : num_dims= 2
MLBosonTagger           INFO    Output0 : dim 0= 1
MLBosonTagger           INFO    Output0 : dim 1= 1
SmoothedContainedWTagg...INFO    SmoothedWZTagger::MLBosonTagger() read value from config
SmoothedContainedWTagg...INFO    SmoothedWZTagger::MLBosonTagger()   + ModelPath /data/BoostedJetTaggers/SmoothedWZTaggers/CNNTaggers/ONNXModel_VVJJ_Model_DisCoJet_13mar22/model.onnx
PathResolver            WARNING Locating dev file dev/MLTest/2020-03-31/t10k-images-idx3-ubyte. Do not let this propagate to a release
PathResolver            WARNING Locating dev file dev/MLTest/2020-03-31/t10k-labels-idx1-ubyte. Do not let this propagate to a release
MLBosonTaggerDisCo      INFO    Using model file: /srv/workDir/usr/WorkDir/0.0.1/InstallArea/x86_64-centos7-gcc8-opt//data/BoostedJetTaggers/SmoothedWZTaggers/CNNTaggers/
ONNXModel_VVJJ_Model_DisCoJet_13mar22/model.onnx
MLBosonTaggerDisCo      INFO    Using pixel file: /cvmfs/atlas.cern.ch/repo/sw/database/GroupData/dev/MLTest/2020-03-31/t10k-images-idx3-ubyte
MLBosonTaggerDisCo      INFO    Using pixel file: /cvmfs/atlas.cern.ch/repo/sw/database/GroupData/dev/MLTest/2020-03-31/t10k-labels-idx1-ubyte
MLBosonTaggerDisCo      INFO    Created the ONNX Runtime session
MLBosonTaggerDisCo      INFO    Output 0 : num_dims= 2
MLBosonTaggerDisCo      INFO    Output0 : dim 0= 1
MLBosonTaggerDisCo      INFO    Output0 : dim 1= 1
SmoothedContainedWTagg...INFO    SmoothedWZTagger::MLBosonTagger() read value from config
SmoothedContainedWTagg...INFO    Smoothed WZ Tagger tool initialized
SmoothedContainedWTagg...INFO      Mass cut low      : (7.180e+01)+(1.124e-02)*pow(x,1)+(-1.721e-05)*pow(x,2)+(4.113e-09)*pow(x,3)
SmoothedContainedWTagg...INFO      Mass cut High     : (1.255e+02)+(-7.448e-02)*pow(x,1)+(5.402e-05)*pow(x,2)+(-1.118e-08)*pow(x,3)
SmoothedContainedWTagg...INFO      DisCoDNN cut low    : (6.731e-01)+(3.434e-05)*pow(x,1)+(1.595e-08)*pow(x,2)+(-7.777e-12)*pow(x,3)
SmoothedContainedWTagg...INFO      Decorate        : 1
SmoothedContainedWTagg...INFO      DecorationName   : SmoothWContained50DisCoJet
SmoothedContainedWTagg...INFO      Pt cut low      : 200
SmoothedContainedWTagg...INFO      Pt cut high     : 3000
SmoothedContainedWTagg...INFO    After tagging, you will have access to the following cuts as a Root::TAccept : (<NCut>) <cut> : <description>)
SmoothedContainedWTagg...INFO      (0)   PassMassLow : mJet > mCutLow
SmoothedContainedWTagg...INFO      (1)   PassMassHigh : mJet < mCutHigh
SmoothedContainedWTagg...INFO      (2)   PassDisCoDNN : DisCoDNNJet > DisCoDNNCut
SmoothedContainedWTagg...INFO    Decorators that will be attached to jet :
SmoothedContainedWTagg...INFO      SmoothWContained50DisCoJet_Cut_mlow : lower mass cut
SmoothedContainedWTagg...INFO      SmoothWContained50DisCoJet_Cut_mhigh : upper mass cut
SmoothedContainedWTagg...INFO    Additional decorators that will be attached to jet :
SmoothedContainedWTagg...INFO      SmoothWContained50DisCoJet_Cut_DisCoDNN : DisCoDNN cut
SmoothedContainedWTagg...INFO    Initializing SmoothedWZTagger tool
```

*Keeping the same messages and debugging as the baseline tagger*

# Onnxruntime Tensor (Ort::Value)

.onnx model takes input in Ort::Value format; (here)
e.g. MNIST hand written digit 3D array (1, 28, 28) needs to be converted to 3D Ort::Value before being fed to a .onnx model with input layer dimension: [-1,28,28]



**Ort::Value::CreateTensor<float>(memory_info, Flatten array , input_tensor_size, [1,28,28], 3)**