# Machine Learning Developments in ROOT

## Sergei Gleyzer, Lorenzo Moneta
### for the ROOT-TMVA Team

**CHEP 2016, October 10, 2016**

# Outline

- **Status and Overview**
- **New TMVA Features**
  - **External Interfaces**
  - **Deep Learning, Jupyter, Parallelization**
- **Future Plans and Outlook**
- **Summary**

# TMVA

**Toolkit for Multivariate Analysis:**

- **HEP Machine Learning workhorse**
- **Part of ROOT**
- **In LHC experiments production**
- **Easy for beginners, powerful for experts**
- **17 active contributors (5 GSoCs)**

# New TMVA version released in upcoming ROOT 6.0.8

# New TMVA Features

# New Features

Modularity, External Interfaces, Updated SVMs
Analyzer Tools: Variable Importance
Deep Learning CPU, GPU
Parallelization with multithreading and GPUs
Analyzer Tools: Cross-Validation,
                        Hyper-Parameter Tuning
Regression Loss Functions
Jupyter: Interactive Training, Visualizations
Unsupervised Learning
Deep Autoencoders
Multi-processing, Spark parallelization
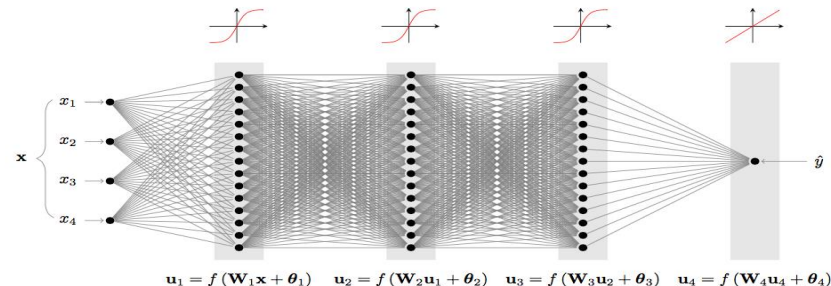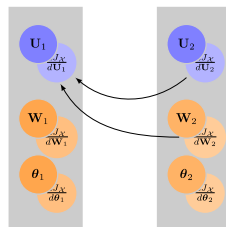
Added in 2015

Added in TMVA ROOT 6.0.8

Upcoming 2016

# TMVA Interfaces

**Interfaces to External ML Tools**

- **RMVA interface to R**

- **PyMVA interface to scikit-learn**

- **KMVA interface to Keras**
  - **High-level interface to Theano, TensorFlow deep-learning libraries**

# Deep Learning

**Is a powerful Machine Learning method based on Deep Neural Networks (DNN) that achieves significant performance improvement in classification tasks**
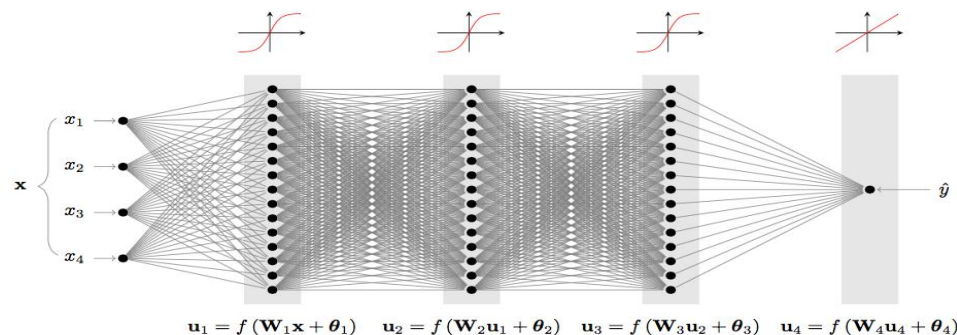


$$\mathbf{u}_1 = f(\mathbf{W}_1 \mathbf{x} + \boldsymbol{\theta}_1) \quad \mathbf{u}_2 = f(\mathbf{W}_2 \mathbf{u}_1 + \boldsymbol{\theta}_2) \quad \mathbf{u}_3 = f(\mathbf{W}_3 \mathbf{u}_2 + \boldsymbol{\theta}_3) \quad \mathbf{u}_4 = f(\mathbf{W}_4 \mathbf{u}_4 + \boldsymbol{\theta}_4)$$

# Deep Learning

## New Deep-Learning Library in TMVA

- **GPU support**
  - CUDA
  - OpenCL



$$\mathbf{u}_1 = f(\mathbf{W}_1\mathbf{x} + \boldsymbol{\theta}_1) \quad \mathbf{u}_2 = f(\mathbf{W}_2\mathbf{u}_1 + \boldsymbol{\theta}_2) \quad \mathbf{u}_3 = f(\mathbf{W}_3\mathbf{u}_2 + \boldsymbol{\theta}_3) \quad \mathbf{u}_4 = f(\mathbf{W}_4\mathbf{u}_4 + \boldsymbol{\theta}_4)$$

- **Excellent performance and high numerical throughput**

# **Deep Learning**

## **CPU Performance:**

### **Implementation:**

- **OpenBLAS, TBB**

### **Peak performance per core:**
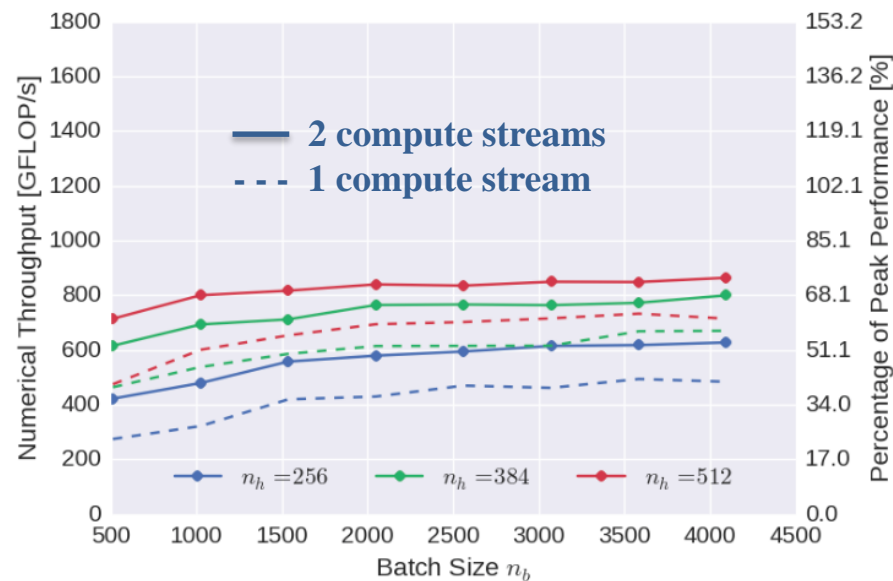
- **16 GFLOP/s**

- **Single, Double Precision**

# **Deep Learning**

**GPU Performance:**

**Network:**

- **20 input nodes**
- **5 hidden layers with $n_h$ nodes each**
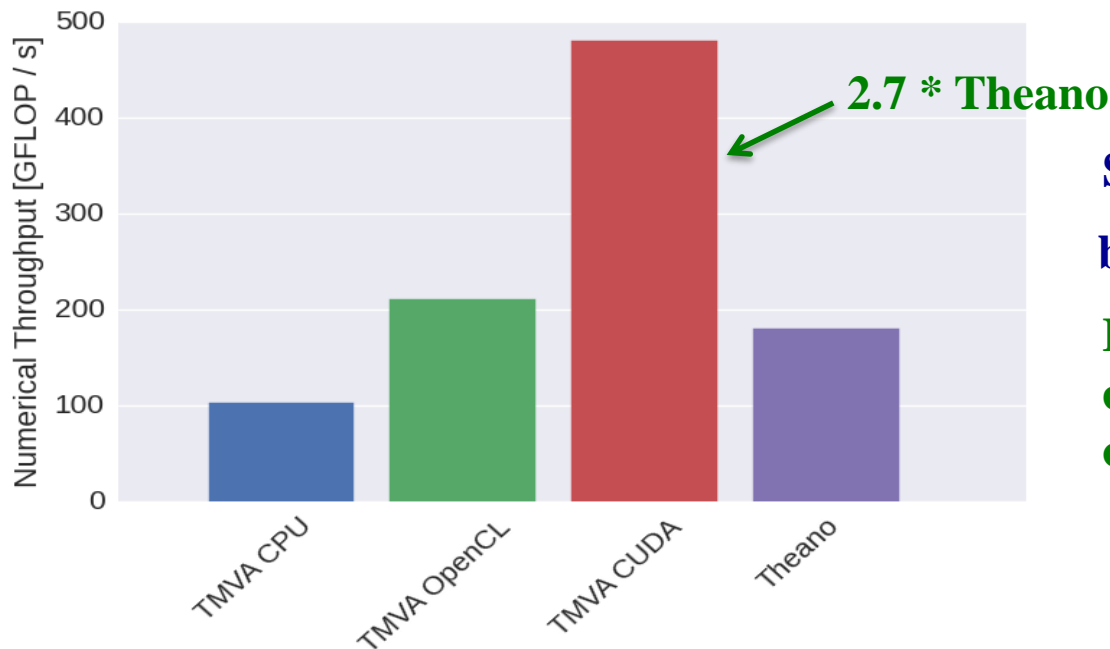
**Hardware:**

- **NVIDIA Tesla K20**
- **1.17 TFLOP/s peak performance @ double precision**



**Good Throughput**

# Deep Learning

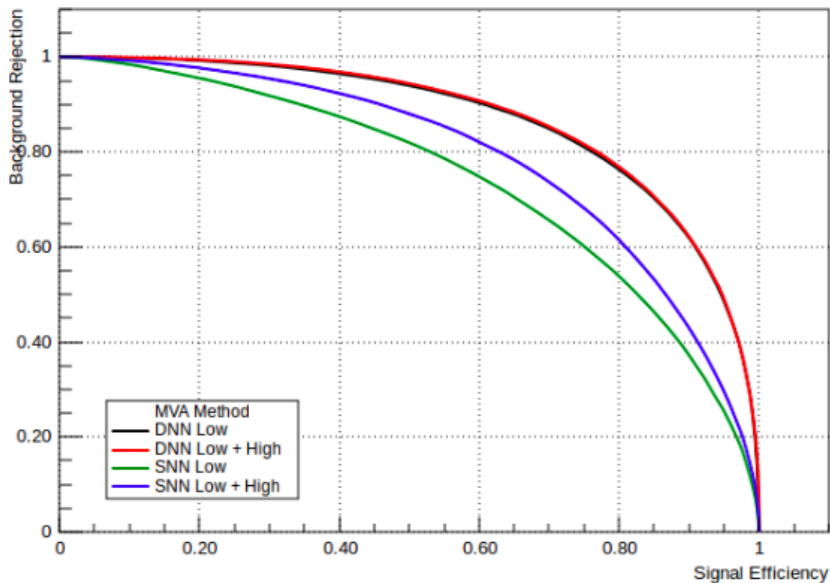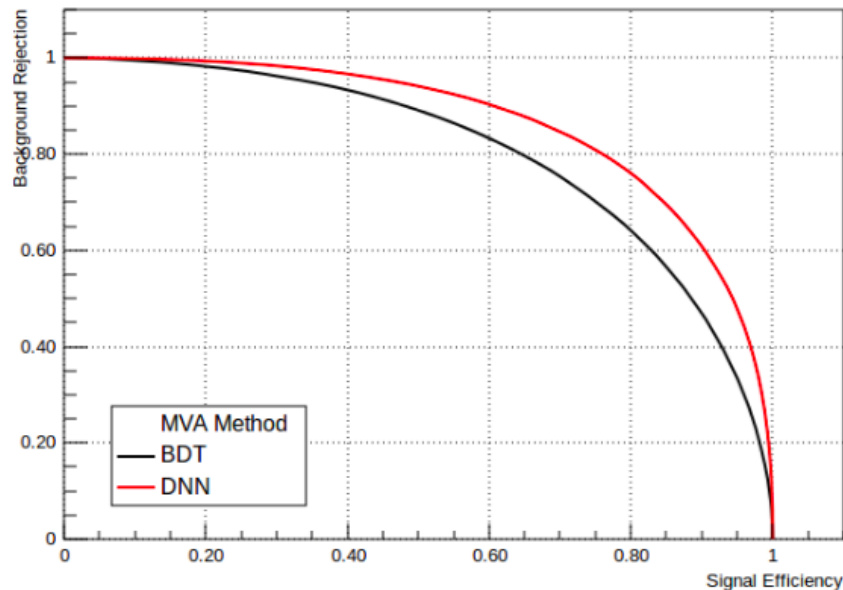## Throughput Comparison



**Single precision**

**batch size = 1024**

**Excellent throughput compared to Theano on same GPU**

# **Deep Learning**



Background Rejection vs. Signal Efficiency
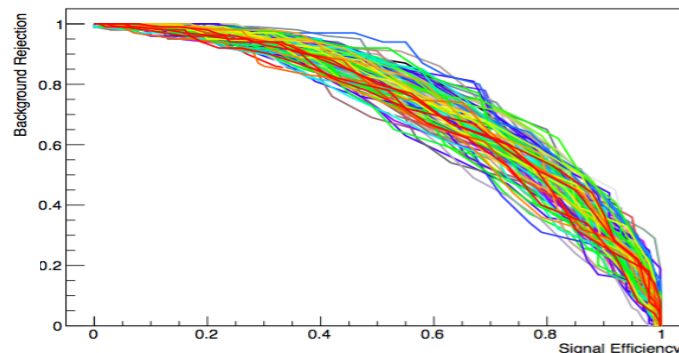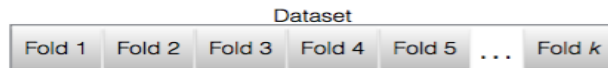
Background Rejection vs. Signal Efficiency

**ROC Performance:** significant improvements compared to shallow networks and boosted decision trees

# Cross Validation

## New features:

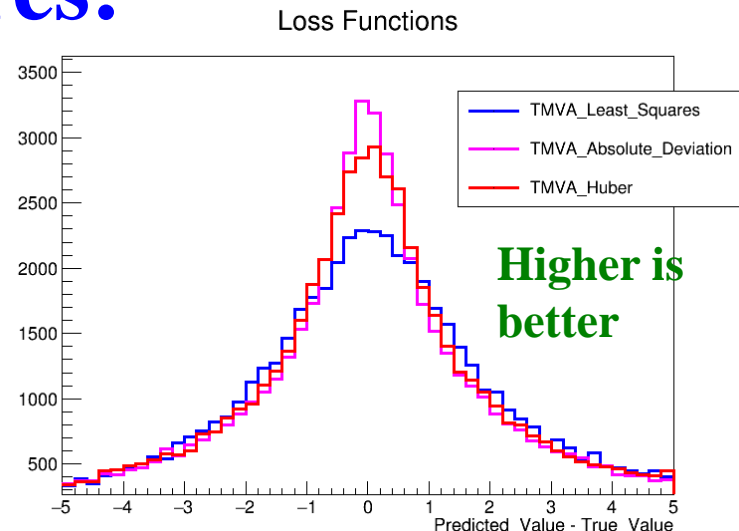- **k-fold cross-validation**



- **Hyper-parameter tuning**
  - **Find optimized parameters (SVM, BDT)**

# Regression

## New Regression Features:

### Loss functions:

- **Huber (default)**
- **Least Squares**
- **Absolute Deviation**
- **Custom Function**

### Important for regression performance



Loss Functions

Higher is better

# **Jupyter Integration**

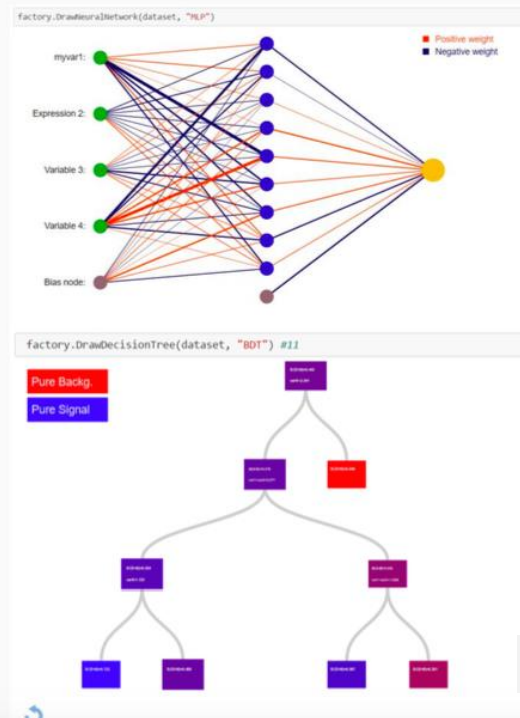## Classifier output: Neural networks, decision trees

### Simple neural network

- Python function reads the network, converts to JSON; JS with d3js make the visualization from JSON
- Interactive: focusing connections, zooming, moving

### Deep neural network

- HTML5 Canvas visualization (speed)
- Less interactive: zooming, moving

### Decision trees

- Ipywidgets: input field for selecting the tree
- Visualization from JSON with D3js
- Interactive: closing subtree, showing the path, focusing, moving, zooming, reset
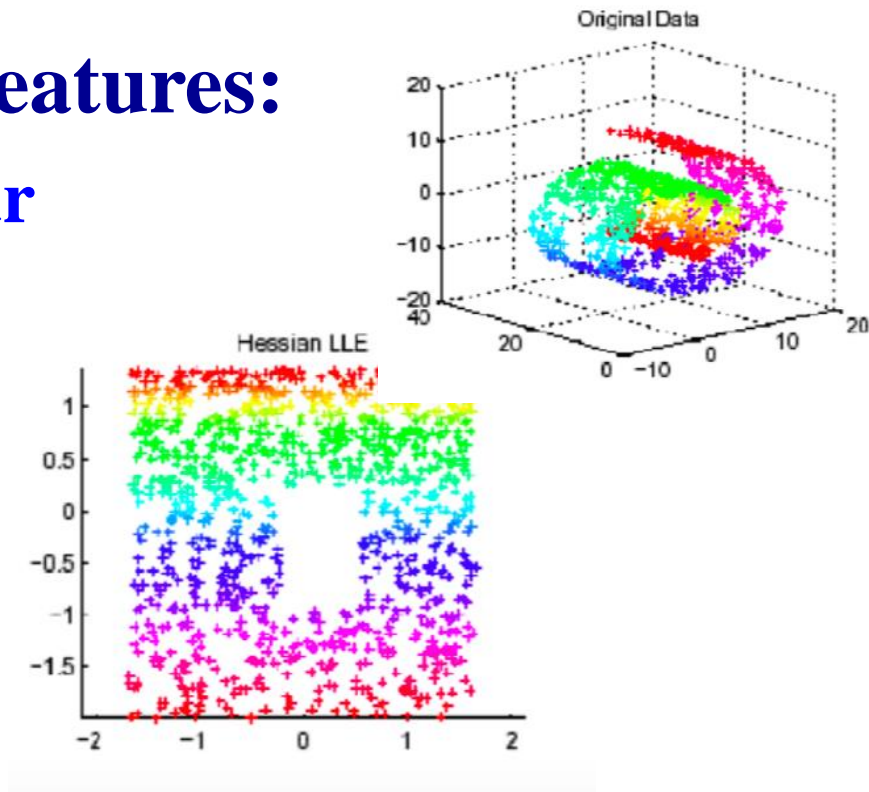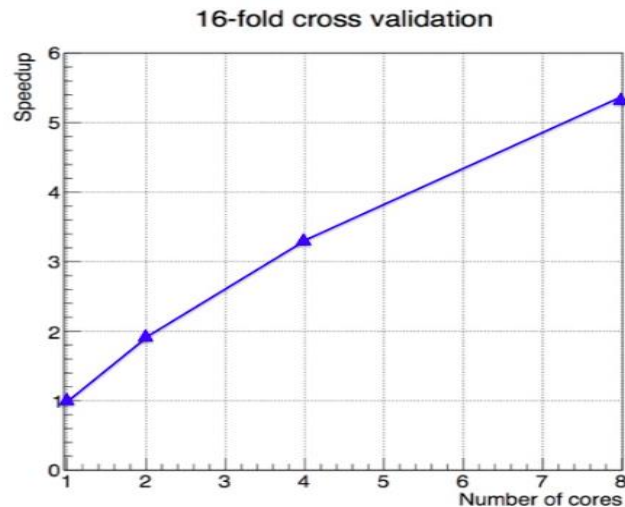
# **Pre-processing**

## **New pre-processing features:**

- **Hessian Locally Linear Embedding**
  - **(Hessian LLE)**
- **Variance Threshold**



Original Data

Hessian LLE

# **Some Upcoming Features**

# **Spark TMVA**

## **SPARK Parallelization**



**Test**

### K-fold cross validation



### 16-fold cross validation



**Good speed-up in prototype R&D**

# Deep Autoencoder



## Deep Autoencoders

x1
x2
x3
x4
x5

Layer 1

a1
a2
a3

Layer 2

$$\begin{bmatrix} a1 \\ a2 \\ a3 \end{bmatrix}$$

New representation for input

- Deep neural network is trained to output the input i.e. learn the identity functions.

- Constrain number of units in hidden layer, thus learning compressed representation.

## Variance Threshold

Input Features → Calculate Variance → Compare Threshold → Selected Features

# **Summary**

- **Many new features in TMVA release upcoming in ROOT 6.0.8**
  - **Production-ready parallelized Deep Learning**
  - **Cross-validation, Hyper-parameter tuning**
  - **Jupyter integration**
  - **More pre-processing features**
  - **Regression updates**
- **Many contributions**
- **Feedback and further contributions welcome**

# Feature Contributors

- **Sergei Gleyzer** — Analyzer Tools, Algorithm Development
- **Lorenzo Moneta** — Multi-threading, Multi-processing
- **Omar Zapata Mesa** — PyMVA, RMVA, Modularity, Parallelization
- **Peter Speckmeyer** — Deep-Learning CPU
- **Simon Pfreundschuh** — Deep-Learning CPU and GPU
- **Adrian Bevan, Tom Stevenson** — SVMs, Cross-Validation, Hyperparameter Tuning
- **Attila Bagoly** — Jupyter Integration, Visualization, Output
- **Albulena Saliji** — TMVA Output Transformation
- **Stefan Wunsch** — KERAS Interfance
- **Pourya Vakilipourtakalou** — Cross-Validation, Parallelization
- **Abhinav Moudhil** — Pre-processing, Deep Autoencoders
- **Georgios Douzas** — Spark, Cross-Validation, Hyperparameter Tuning
- **Paul Seyfert** — Performance optimization of MLP
- **Andrew Carnes** — Regression, Loss Functions, BDT Parallelization

**Continued invaluable contributions from Andreas Hoecker, Helge Voss, Eckhard von Thorne, Jörg Stelzer, and key support from CERN EP-SFT Group**

# **More Information**

**Websites:** **http://root.cern.ch**
**http://iml.cern.ch**
**http://oproject.org**
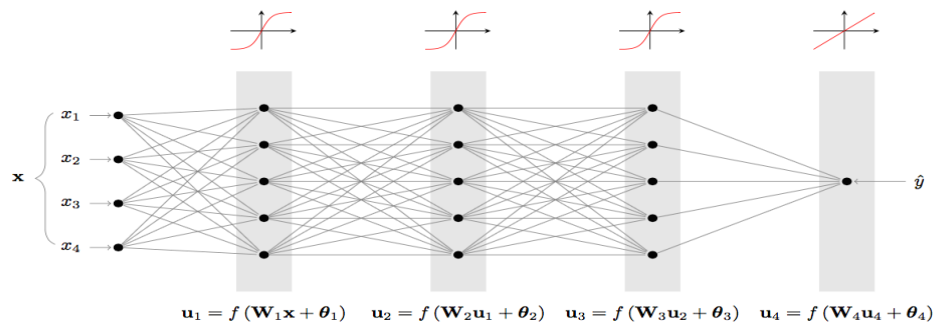
# Inter-experimental LHC Machine Learning working group

- **Exchange of HEP-ML expertise and experience among LHC experiments**
- **ML Forum**
- **ML software development and maintenance**
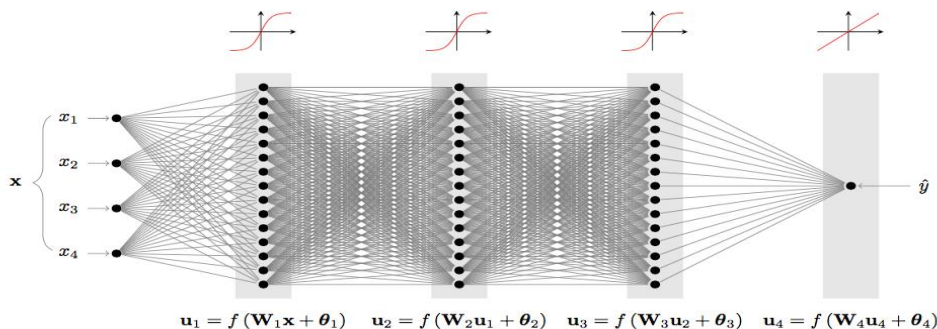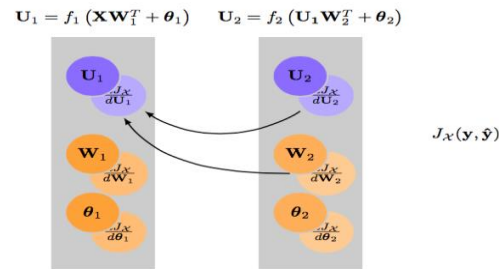- **Exchange between HEP and ML communities**
- **Education (Tutorials)**

# **Backup**

Sergei V. Gleyzer CHEP 2016

# TMVA Deep Learning



$$\mathbf{u}_1 = f\left(\mathbf{W}_1\mathbf{x} + \boldsymbol{\theta}_1\right) \quad \mathbf{u}_2 = f\left(\mathbf{W}_2\mathbf{u}_1 + \boldsymbol{\theta}_2\right) \quad \mathbf{u}_3 = f\left(\mathbf{W}_3\mathbf{u}_2 + \boldsymbol{\theta}_3\right) \quad \mathbf{u}_4 = f\left(\mathbf{W}_4\mathbf{u}_4 + \boldsymbol{\theta}_4\right)$$

**Design**

$$\mathbf{U}_1 = f_1\left(\mathbf{X}\mathbf{W}_1^T + \boldsymbol{\theta}_1\right) \quad \mathbf{U}_2 = f_2\left(\mathbf{U}_1\mathbf{W}_2^T + \boldsymbol{\theta}_2\right)$$

$$\begin{bmatrix} x_{0,0} & \cdots & x_{0,m} \\ x_{1,0} & \cdots & x_{1,m} \\ \vdots & & \vdots \\ x_{n,0} & \cdots & x_{n,m} \end{bmatrix}$$

$$J_{\mathcal{X}}(\mathbf{y}, \hat{\mathbf{y}})$$

link

$$\mathbf{u}_1 = f\left(\mathbf{W}_1\mathbf{x} + \boldsymbol{\theta}_1\right) \quad \mathbf{u}_2 = f\left(\mathbf{W}_2\mathbf{u}_1 + \boldsymbol{\theta}_2\right) \quad \mathbf{u}_3 = f\left(\mathbf{W}_3\mathbf{u}_2 + \boldsymbol{\theta}_3\right) \quad \mathbf{u}_4 = f\left(\mathbf{W}_4\mathbf{u}_4 + \boldsymbol{\theta}_4\right)$$