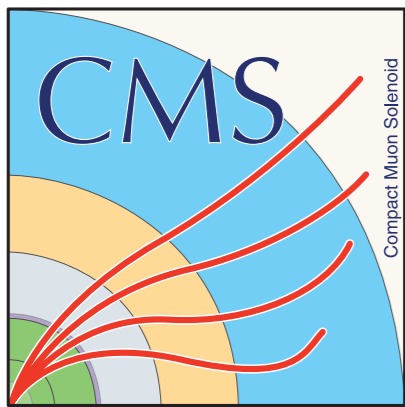# MACHINE LEARNING INFERENCE IN CMSSW

Huilin Qu

*on behalf of the CMS collaboration*

*ATLAS Machine Learning Workshop*
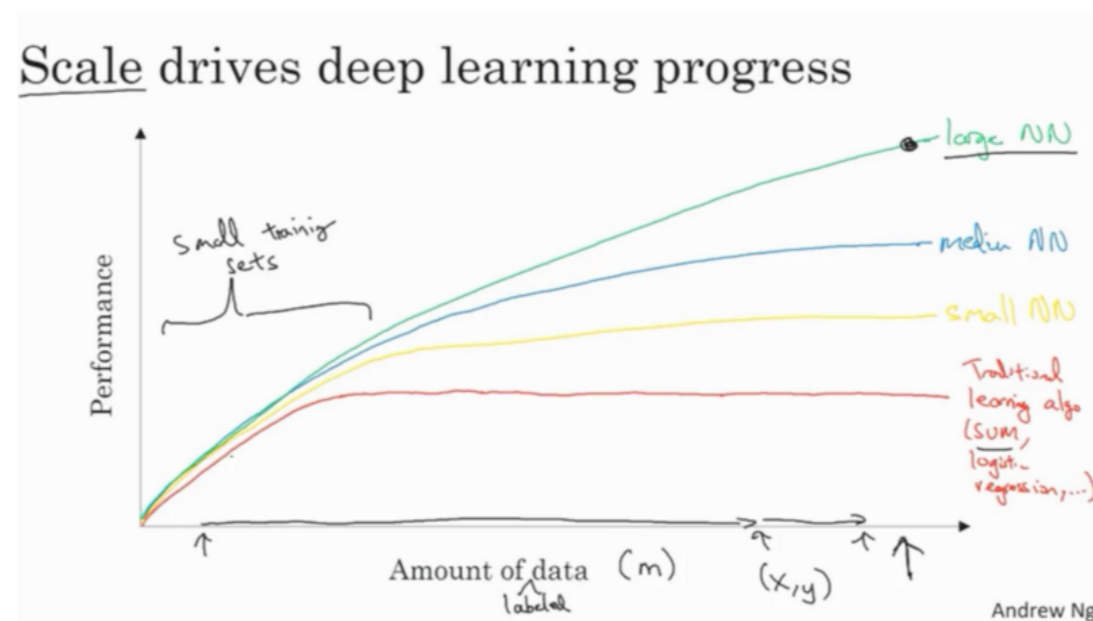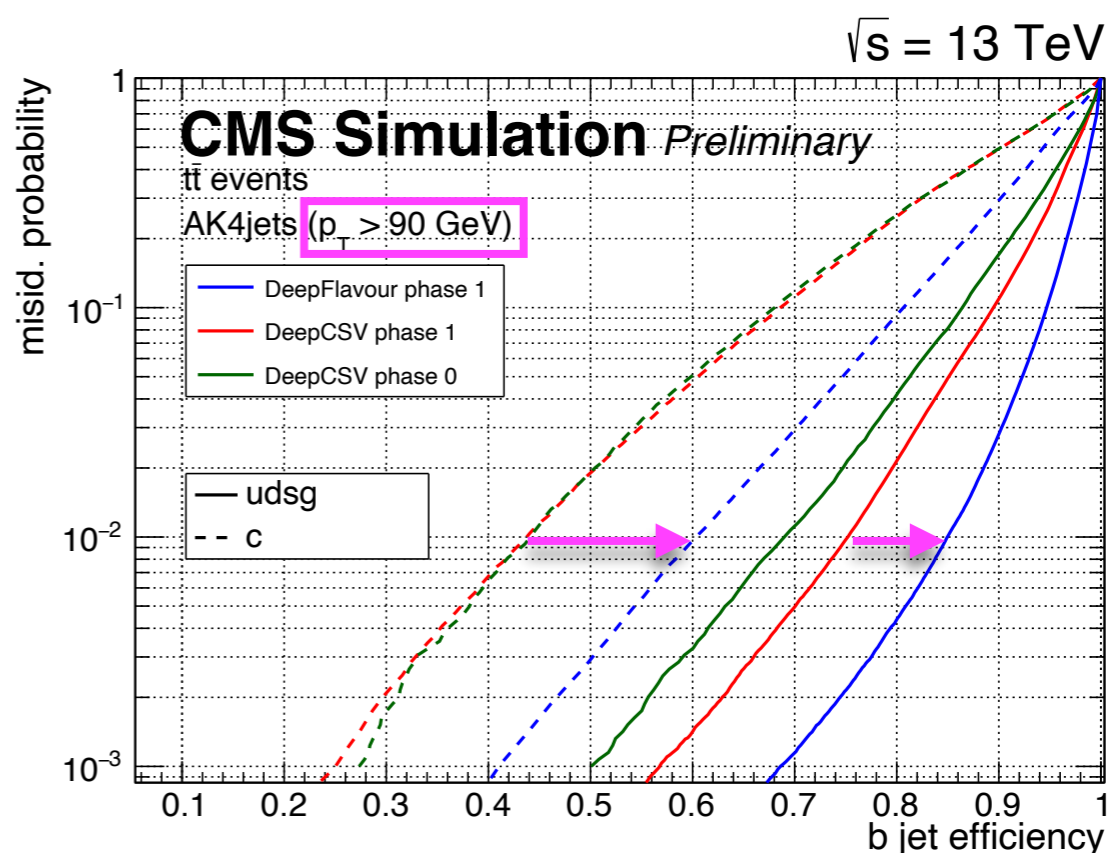
*November 15, 2019*

# INTRODUCTION

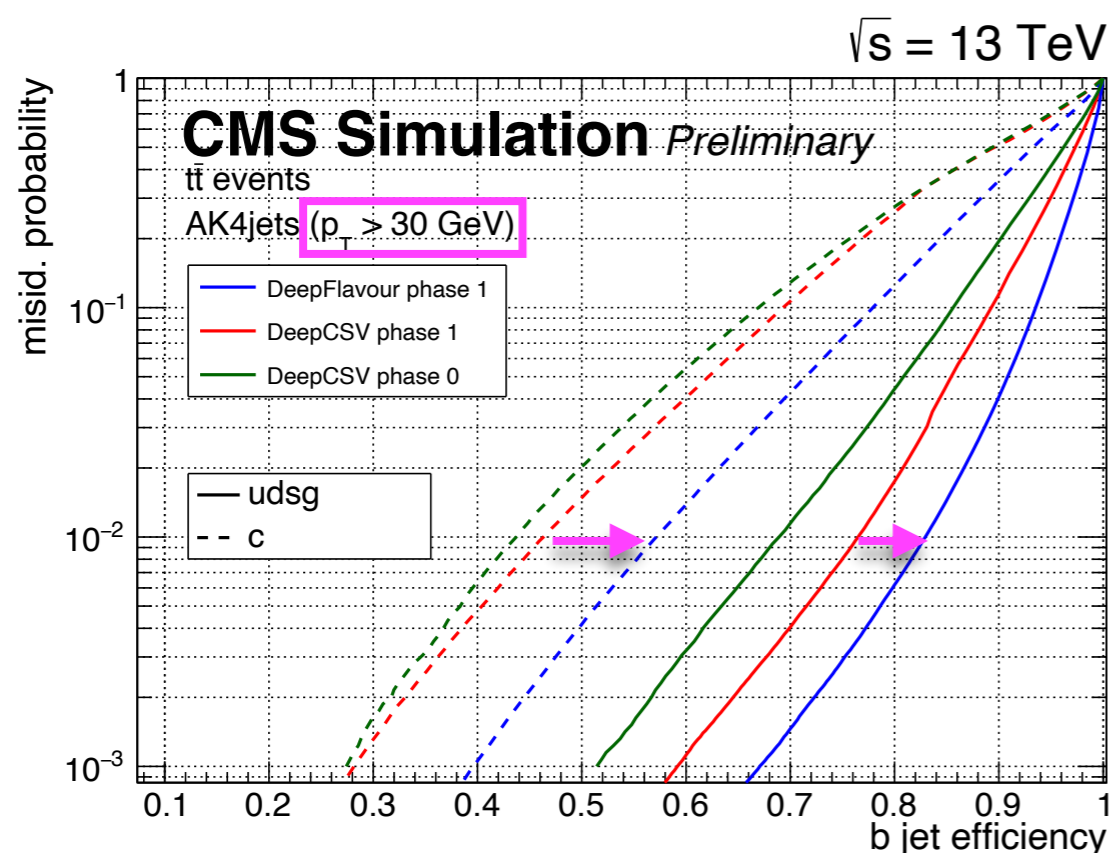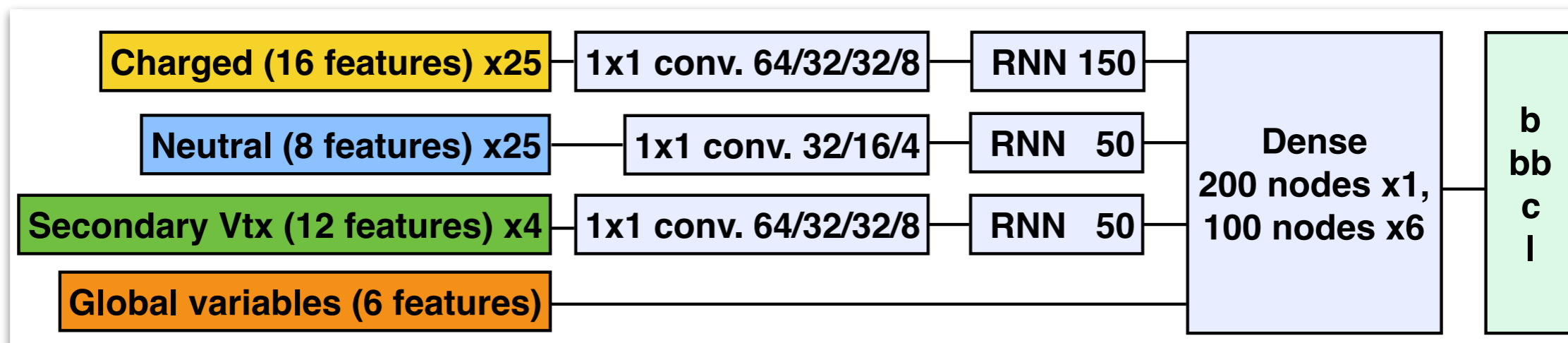- Machine learning (ML) can provide powerful tools for particle physics experiments

- Trend in recent years: deep learning (DL) + low-level inputs



Scale drives deep learning progress

- A variety of new DL algorithms have been developed in CMS

  - b-tagging

  - boosted jet tagging
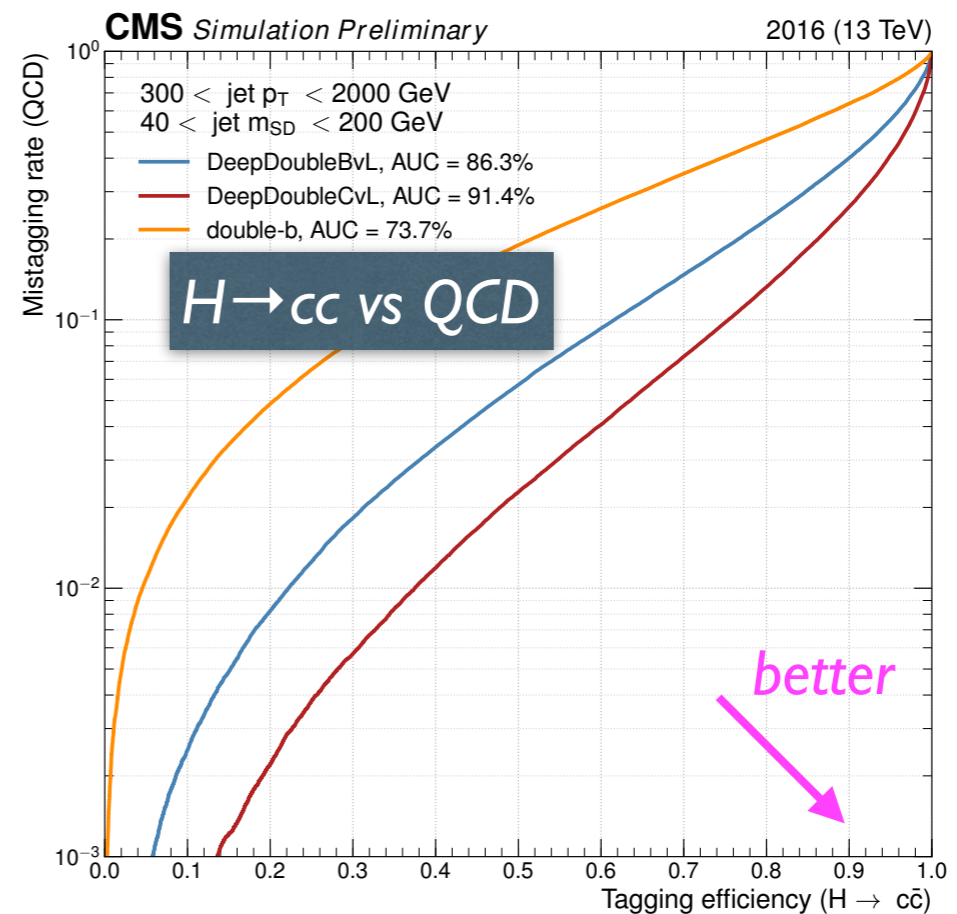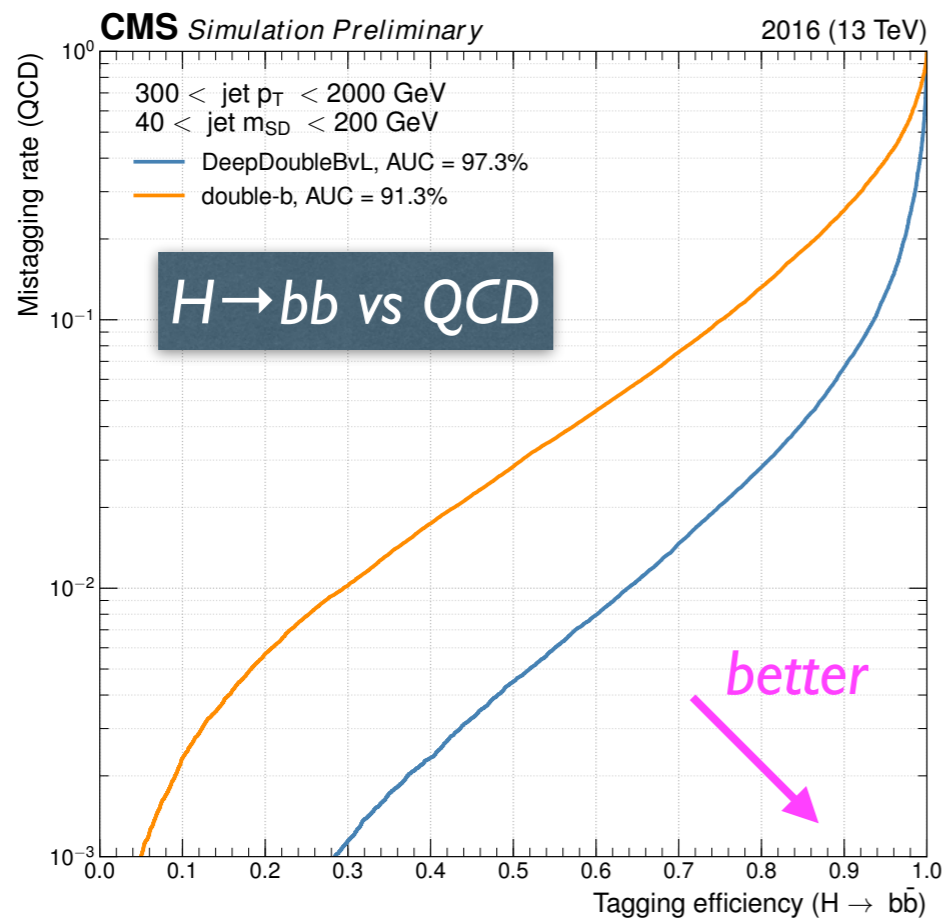
  - tau identification

  - b-jet energy regression

  - …

# DEEPJET (DEEPFLAVOUR)

- AK4 jet flavour tagger

CMS-DP-2018-033

3

# DEEPDOUBLEX

- Boosted jet flavour tagger for bb/cc

**CMS** *Simulation Preliminary*                     2016 (13 TeV)

$300 < $ jet $p_T < 2000$ GeV
$40 < $ jet $m_{SD} < 200$ GeV

— DeepDoubleBvL, AUC = 97.3%
— double-b, AUC = 91.3%

*H→bb vs QCD*

*better*

Mistagging rate (QCD)

Tagging efficiency (H → b$\bar{b}$)

Mass Sculpting
DeepDoubleBvL

**Figure 4.** Effect on the jet soft-drop mass distribution of misidentified events by the DeepDoubleBvL identification algorithm demonstrating the degree to which the algorithm is dependent on the mass of the jet. These histograms are obtained for a fixed overall mistagging rate from a QCD sample.

7/5/2018

**CMS** *Simulation Preliminary*                     2016 (13 TeV)

$300 < $ jet $p_T < 2000$ GeV
$40 < $ jet $m_{SD} < 200$ GeV

— DeepDoubleBvL, AUC = 86.3%
— DeepDoubleCvL, AUC = 91.4%
— double-b, AUC = 73.7%

*H→cc vs QCD*

*better*

Mistagging rate (QCD)

Tagging efficiency (H → c$\bar{c}$)

*CMS-DP-2018-046*

10

4

# DEEPAK8

- ## Multi-class boosted jet tagger for top / W / Z / H

### Inputs

**Substructure**

**Particles**
- Up to 100 PF candidates[(*)]
- Sorted in descending $p_T$ order
- Uses basic kinematic variables, Puppi weights, and track properties (quality, covariance, displacement, etc.)

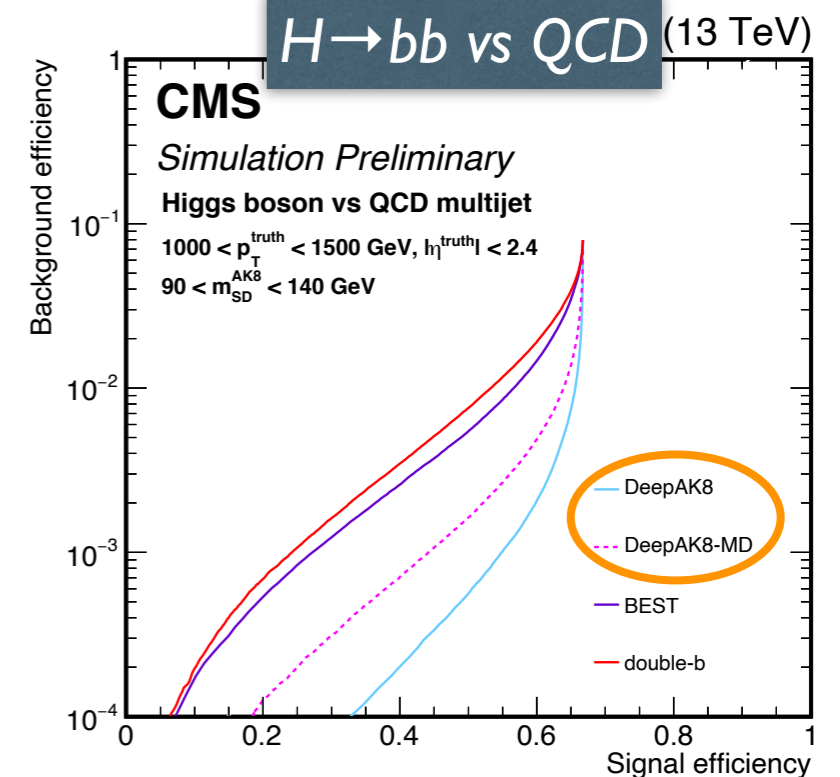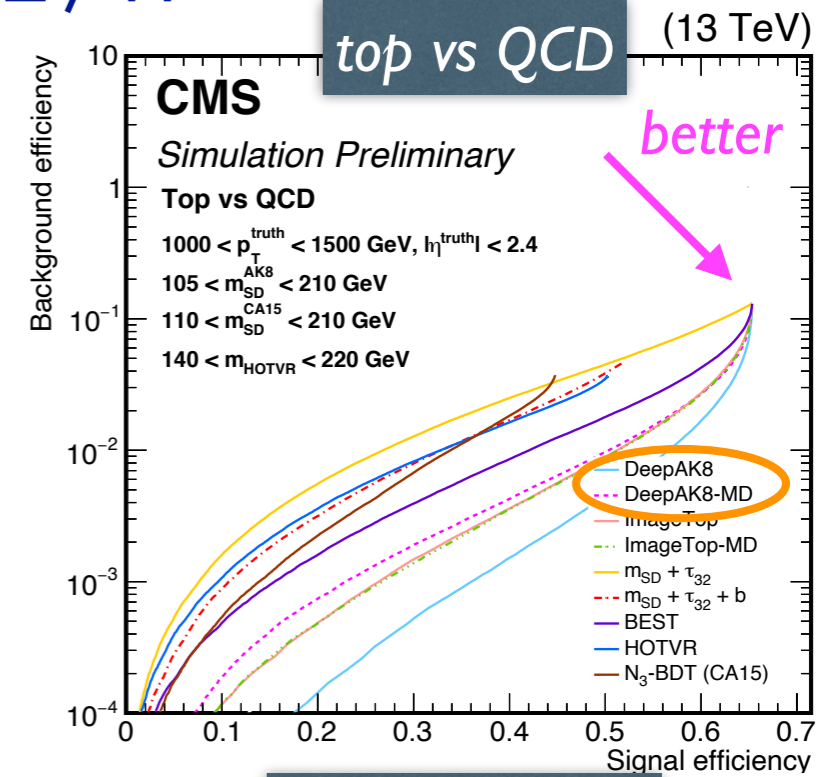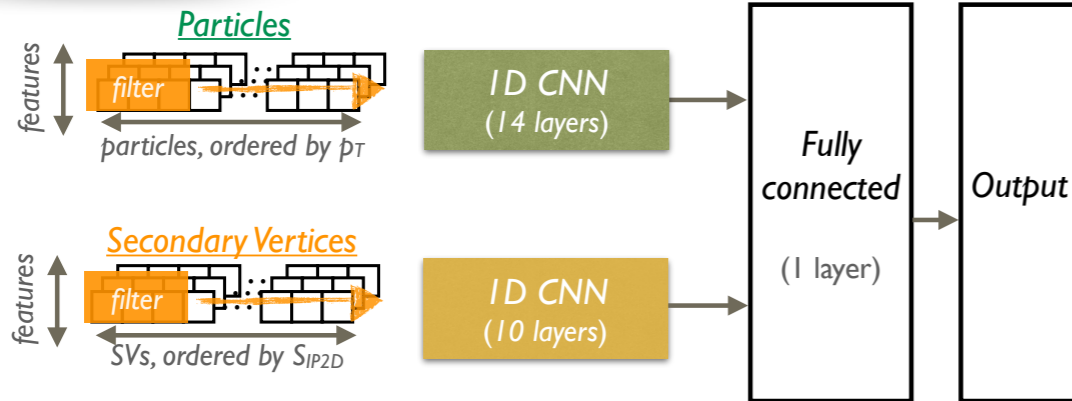**Flavour**

**Secondary vertices**
- Up to 7 SVs[(*)] (inside jet cone)
- Sorted in descending $S_{IP2D}$ order
- Uses SV kinematics and properties (quality, displacement, etc.)

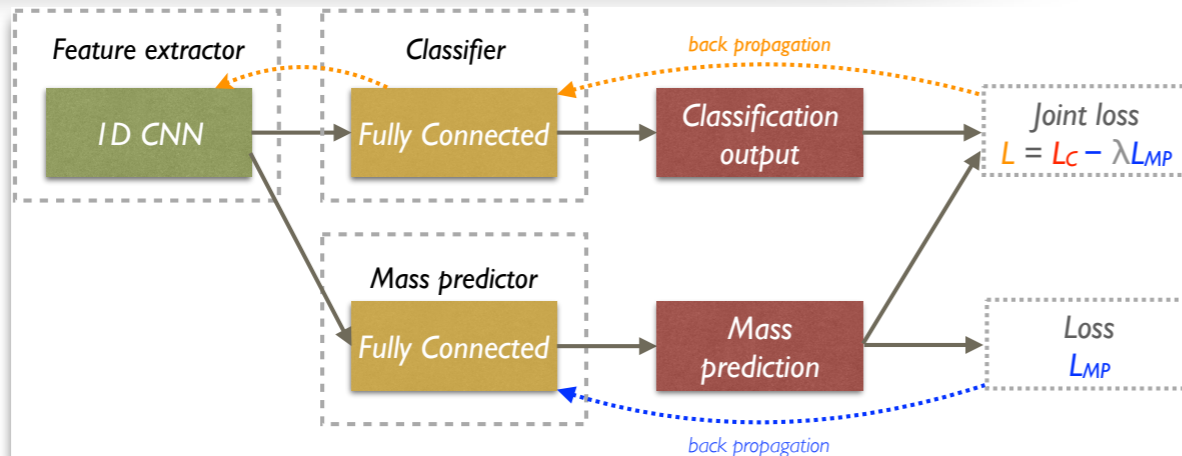*(*) Number chosen to include all candidates for ≥ 90% of the events*

### Architecture

**Particles**

features → filter → 1D CNN (14 layers) → 

particles, ordered by $p_T$

**Secondary Vertices**

features → filter → 1D CNN (10 layers) → 

SVs, ordered by $S_{IP2D}$

→ Fully connected (1 layer) → Output

### Output

| Category | Label |
|----------|-------|
| | H (bb) |
| Higgs | H (cc) |
| | H (VV*→qqqq) |
| | top (bcq) |
| Top | top (bqq) |
| | top (bc) |
| | top (bq) |
| W | W (cq) |
| | W (qq) |
| | Z (bb) |
| Z | Z (cc) |
| | Z (qq) |
| | QCD (bb) |
| | QCD (cc) |
| QCD | QCD (b) |
| | QCD (c) |
| | QCD (others) |

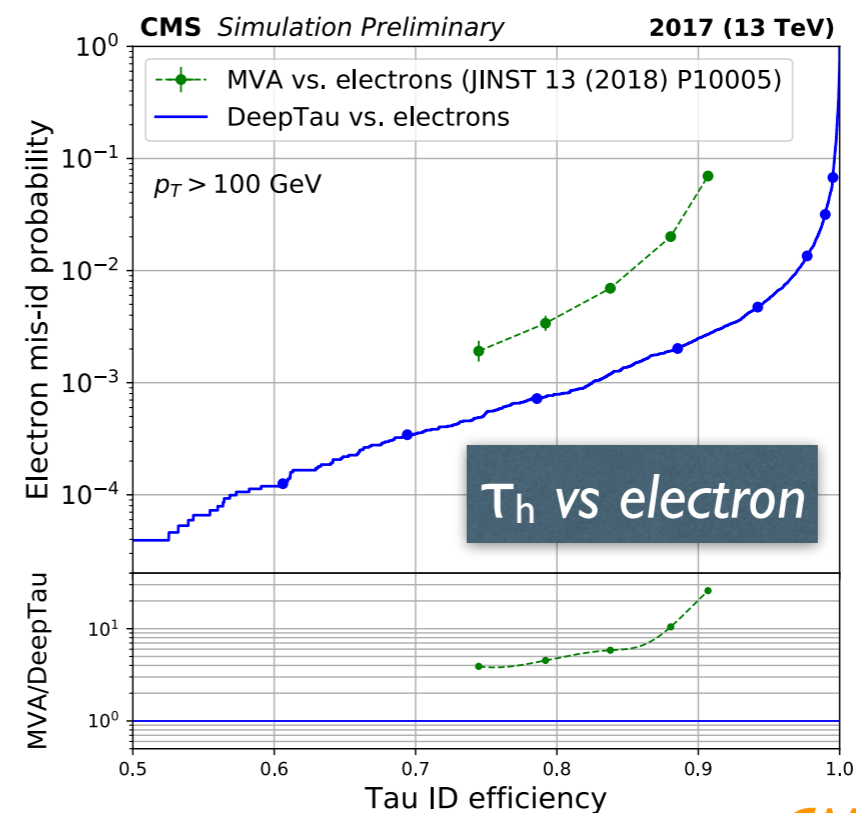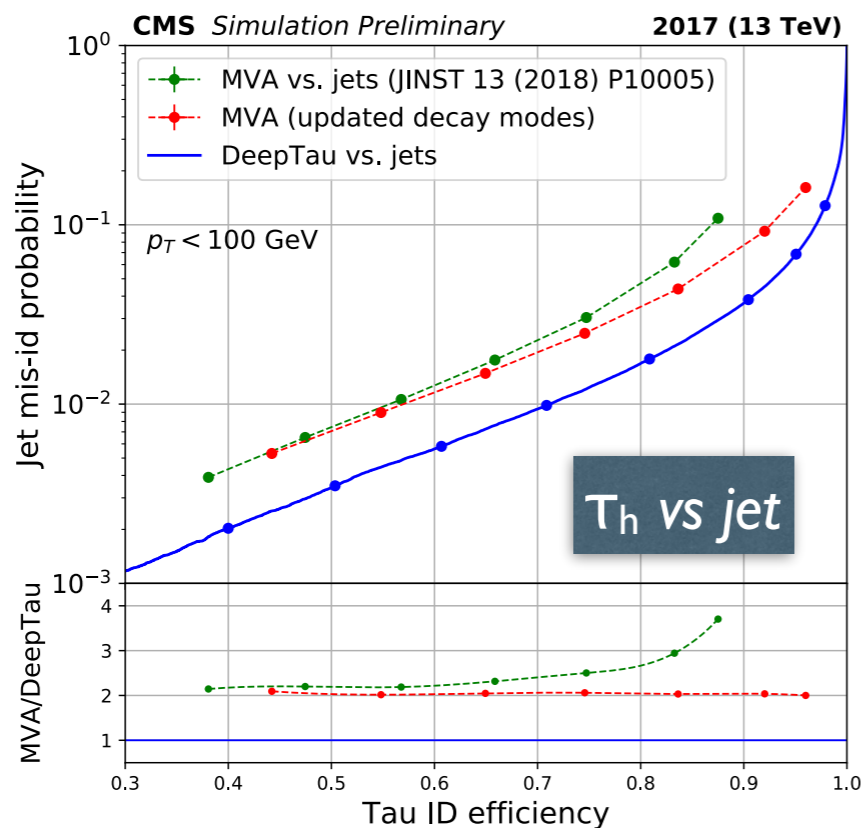*DeepAK8-MD: mass decorrelation w/ adversarial training*

**Feature extractor**

1D CNN

**Classifier**

Fully Connected → Classification output

back propagation

Joint loss $L = L_C - \lambda L_{MP}$

**Mass predictor**

Fully Connected → Mass prediction

Loss $L_{MP}$

back propagation



**top vs QCD** (13 TeV)

*better*

**CMS**
*Simulation Preliminary*
**Top vs QCD**
$1000 < p_T^{truth} < 1500$ GeV, $|\eta^{truth}| < 2.4$
$105 < m_{SD}^{AK8} < 210$ GeV
$110 < m_{SD}^{CA15} < 210$ GeV
$140 < m_{HOTVR} < 220$ GeV

- DeepAK8
- DeepAK8-MD
- ImageTop
- ImageTop-MD
- $m_{SD} + \tau_{32}$
- $m_{SD} + \tau_{32} + b$
- BEST
- HOTVR
- $N_3$-BDT (CA15)

Background efficiency / Signal efficiency

**H→bb vs QCD** (13 TeV)

**CMS**
*Simulation Preliminary*
**Higgs boson vs QCD multijet**
$1000 < p_T^{truth} < 1500$ GeV, $|\eta^{truth}| < 2.4$
$90 < m_{SD}^{AK8} < 140$ GeV

- DeepAK8
- DeepAK8-MD
- BEST
- double-b

Background efficiency / Signal efficiency

CMS-DP-2017-049   CMS-PAS-JME-18-002   5

# DEEPTAU

- CNN-based hadronic tau identification algorithm

**Input cells**

**hadrons block**
38 variables per cell

**e-gamma block**
86 variables per cell

**muon block**
64 variables per cell

**Inner cone = signal cone**
$$\Delta R < \max\left[\min\left(0.1, \frac{3}{p_T^\tau}\right), 0.05\right]$$

**Outer cone = signal & isolation cones**
$$\Delta R < 0.5$$

Total number of trainable parameters (TP) in the DeepTau network = **1 155 353**

| | *208 819 TP* | *185 600 TP* | *161 004 TP* |

Inner cells → Pre-processing each inner cell separately → 5 convolution of inner cells with 3×3 windows →

| | *208 819 TP* | *371 200 TP* | |

Outer cells → Pre-processing each outer cell separately → 10 convolutions of outer cells with 3×3 windows → 5 dense layers → $\begin{matrix} p_e \\ p_\mu \\ p_\tau \\ p_j \end{matrix}$

High level features → Pre-processing of high level features →

*19 911 TP*

**CMS** *Simulation Preliminary*          **2017 (13 TeV)**

- ● MVA vs. electrons (JINST 13 (2018) P10005)
- — DeepTau vs. electrons

$p_T < 100$ GeV

Electron mis-id probability

τ$_h$ *vs jet*

MVA/DeepTau

Tau ID efficiency

**CMS** *Simulation Preliminary*          **2017 (13 TeV)**

- ● MVA vs. electrons (JINST 13 (2018) P10005)
- — DeepTau vs. electrons

$p_T > 100$ GeV

Electron mis-id probability

τ$_h$ *vs electron*

MVA/DeepTau

Tau ID efficiency

# B-JET ENERGY REGRESSION

- Simultaneous estimation of the b-jet energy and its resolution

43 Inputs

**3 Outputs**

$\hat{y}$   - correction

$\hat{y}_{25\%}$ - 25% quantile

$\hat{y}_{75\%}$ - 75% quantile

resolution estimator $\hat{\sigma} = \dfrac{\hat{y}_{75\%} - \hat{y}_{25\%}}{2}$

*More details in*
*N. Chernyavskaya's talk*

**Joint loss function for correction (Huber) and resolution (quantiles) :**

$$Loss = Huber(y, F(x)) + \rho_{0.75}(y - F(x)) + \rho_{0.25}(y - F(x))$$

**(13 TeV)**

Entries / 5 GeV

**CMS** *Simulation Preliminary*

- DNN
  - μ = 124.6 GeV
  - σ = 15.4 GeV
- Baseline
  - μ = 115.9 GeV
  - σ = 18.0 GeV

$m_{jj}$ (GeV)

*13% improvement in per-jet relative resolution*
*20% improvement in dijet mass resolution*

*Successfully applied to the CMS*
*H→bb observation analysis*

*CMS-PAS-HIG-18-027*

7

# FROM DEVELOPMENT TO DEPLOYMENT

- The development of a DL model takes lots of effort

  - *a good DL model = input feature selection + training dataset preparation + network architecture design + hyperparameter optimization + ...*

- Next step: deploying to production!

  - but...

# FROM DEVELOPMENT TO DEPLOYMENT

- The development of a DL model takes lots of effort

  - *a good DL model = input feature selection + training dataset preparation + network architecture design + hyperparameter optimization + …*

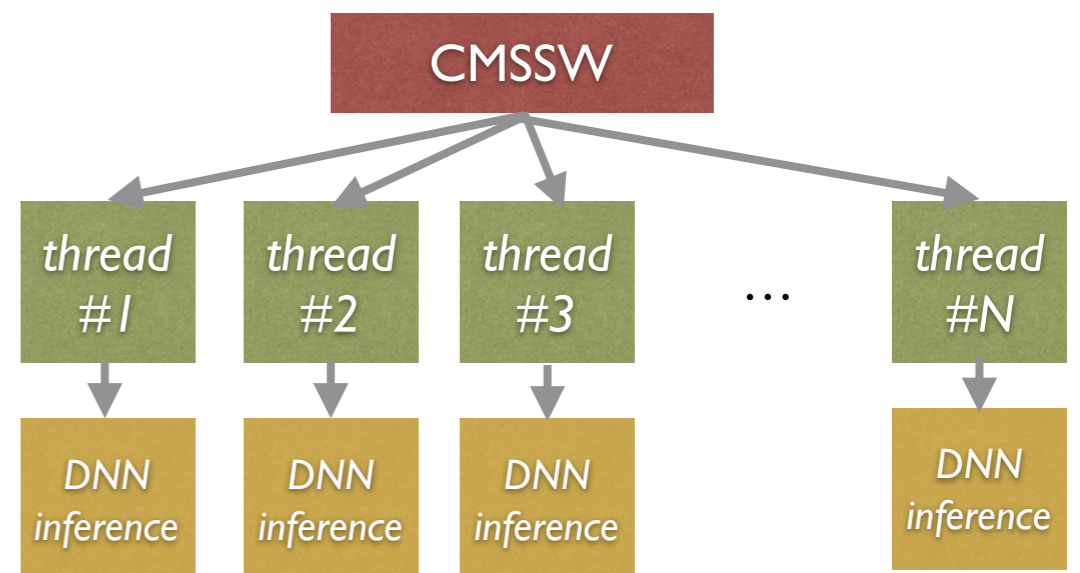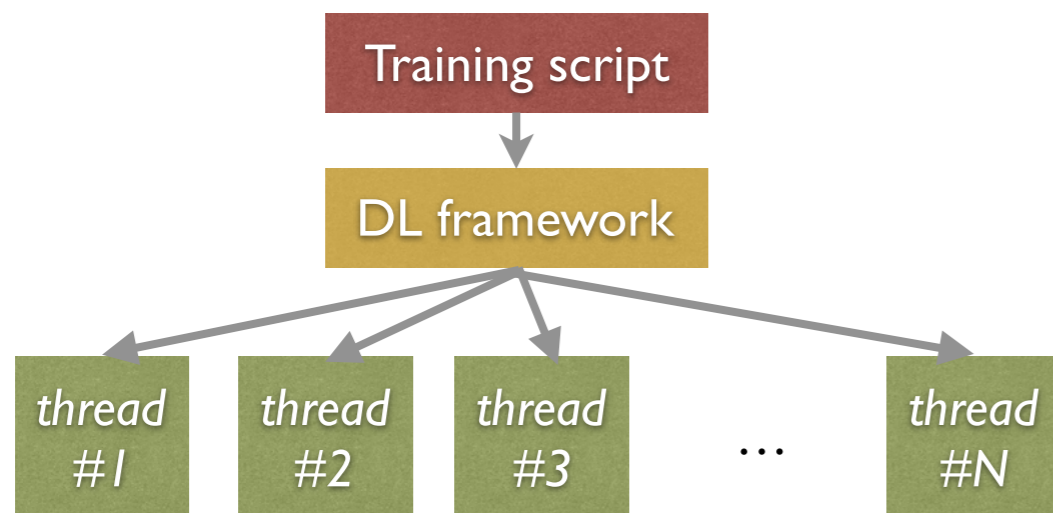- Next step: deploying to production!

  - but…

# TRAINING VS INFERENCE

- Training

  - python based

  - typically on GPU

  - large batch size

    - O(100) to O(1000)

  - standalone environment

    - weak constraints on timing/ memory

    - multi-threading managed by the DL framework (e.g., TensorFlow)

- Inference

  - C/C++ based

  - on CPU

  - small batch size

    - O(1) to O(10), often just 1

  - integrated in the experimental software (e.g., CMSSW)

    - tight constraints on timing/memory

    - multi-threading managed by the experimental software

# ML INFERENCE ENGINES IN CMSSW

- A number of DL frameworks have been integrated into CMSSW to support the new DNN-based algorithms

  - TensorFlow

    - DeepJet

    - DeepDoubleX

    - DeepTau

    - b-jet energy regression

  - MXNet

    - DeepAK8

  - ONNX Runtime

    - new development

    - can support all these models by converting to ONNX format

- Modifications are needed for all of them to work nicely with CMSSW

# TensorFlow

- TensorFlow

    - most widely used framework (together with Keras)

    - complicated to build and integrate (requires Google's build system bazel, etc.)

- Issue 1: Multi-Threading

    - upon startup, TF creates *lots* of threads in its thread pool for parallel data loading and parallelism within/between operators

    - good for end-users who runs only 1 thread to call TF, but not good for HEP frameworks that typically manage their own threading schemes (CMSSW uses TBB)

    - solved with the implementation of two custom sessions

        - NTSession (default in CMSSW): disable multi-threading

        - TBBSession: threads scheduled by Intel's TBB

- Issue 2: Memory Consumption

    - TF Graphs obtained after training can be quite large (e.g., 150 MB for DeepJet)

    - memory footprint can be reduced by a factor of O(10-100) by:

        - converting *variables* to *constant* tensors (`freeze_graph`)

        - removing ancillary information needed only for training

        - a number of tools available online, e.g., keras_to_tensorflow

    - further reduction: load TF graph only once and share it among all threads (sessions)

# MXNᴇᴛ

- MXNet

    - a DL framework focused on efficiency and scalability

    - relatively straightforward to build and integrate (cmake build system, standard BLAS library)

    - exported models are ready-to-use for inference (model json + binary parameter file)

- Issue 1: multi-threading

    - similar problem as TF: MXNet creates and manages its own thread pool

    - <u>solution</u>: use MXNet's "NaiveEngine" (no threading) and make it "<u>thread_local</u>" (so each thread can call it independently)

        - need to <u>re-assign the resources</u> (workspace) in each run to ensure thread-safety (more details in <u>M. Verzetti's talk</u> last year)

- Issue 2: BLAS library

    - DeepAK8 inference runs 4-5x slower in CMSSW than w/ standalone MXNet

    - the problem was tracked to the use of the BLAS library

        - the standalone MXNet links to OpenBLAS statically

        - MXNet in CMSSW is built to link with OpenBLAS dynamically, but a slower BLAS library (glsblas) is used by other softwares (e.g., ROOT) and loaded first, thus providing the BLAS symbols to MXNet

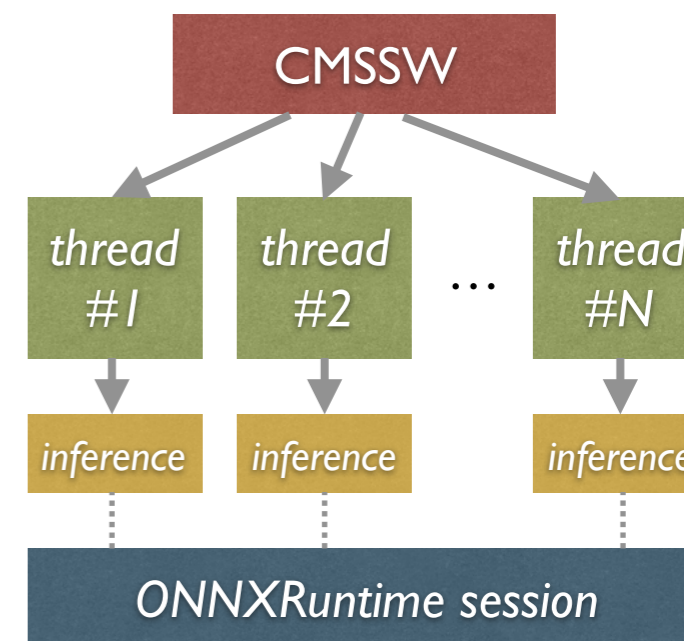        - solved by linking to OpenBLAS consistently in all CMS softwares

# ONNX RUNTIME

- Open Neural Network Exchange (ONNX)

    - an open source format for ML models w/ increasing adoption

    - supports most of the main-stream DL operators

    - conversion tools available for most of the DL frameworks: Keras/TF, PyTorch, MXNet, etc.

- ONNX Runtime

    - "a performance-focused complete scoring engine for ONNX models"

    - advantages:

        - flexibility: can support a wide range of models via ONNX

        - speed: optimized for inference (including on CPUs), rather than training (TF/MXNet/PyTorch etc.)

        - thread-safety: "Multiple threads can invoke the Run() method on the same inference session object."

    - caveats:

        - ONNX may not support all, especially novel ML models

# ONNX RUNTIME INTEGRATION

- **A few modifications to make it work better w/ CMSSW**

  - configured it to run in a "no-threading" mode

    - i.e., each CMSSW thread uses the global inference session object to run inference concurrently with no extra threads

    - setting `intra_op_num_threads` and `inter_op_num_threads` to 1 gives the desired behavior (i.e., it does not create any new threads)

    - however, need to <u>remove a hard-coded thread pool</u> for the LSTM operator

      - likely will be fixed officially in the future

  - introduced an environment variable to control the runtime kernel selection

    - ONNX Runtime's math library (MLAS) selects the fastest compute kernel *dynamically* based on the available CPU instruction sets

      - outputs are not bitwise equal on different CPU architectures as different instructions (SSE/AVX/AVX2/etc.) will be used — causes trouble for PR validation

    - <u>added an environment variable</u> to control the allowed instruction sets

      - default to using only SSE: not attempting to use more advanced instructions (like AVX)

      - ensures bitwise reproducibility across different CPU architectures

      - dynamic kernel selection can be switched on for production to save run time
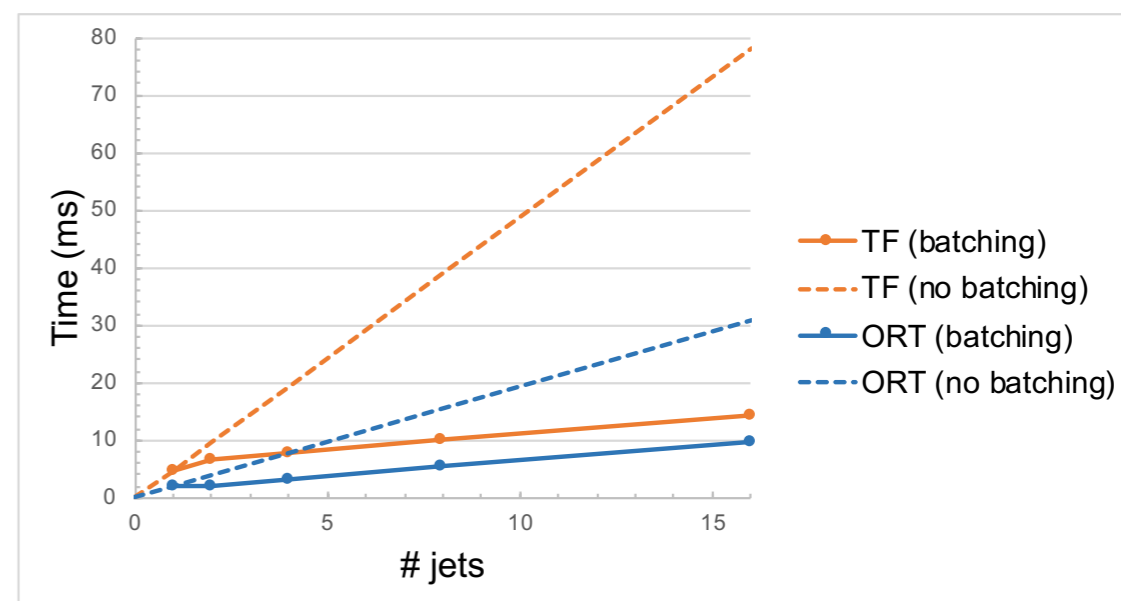
# ONNX RUNTIME: TIMING

- Significant speed-up w/ ONNX Runtime compared to the TF/MXNet based implementation ([cms-sw/cmssw#28112](cms-sw/cmssw#28112))

  - depending on the network architecture, the speed-up varies from 10-15% to ~10x

  - the use of newer vector instructions (e.g., AVX) can bring further improvements

| Time (s) / event | Baseline | ONNX Runtime (SSE) | Speed-up w.r.t baseline | ONNX Runtime (AVX) | Speed-up w.r.t baseline |
|---|---|---|---|---|---|
| DeepTau | 0.039245 | 0.053057 | 0.74 | 0.024901 | 1.58 |
| DeepJet | 0.058576 | 0.009333 | 6.28 | 0.006735 | 8.70 |
| DeepAK8 | 0.003538 | 0.003222 | 1.10 | 0.002107 | 1.68 |
| DeepAK8-MD | 0.003598 | 0.003153 | 1.14 | 0.002078 | 1.73 |
| DeepDoubleBvL | 0.004457 | 0.000451 | 9.88 | 0.000363 | 12.28 |
| DeepDoubleCvB | 0.004514 | 0.000445 | 10.14 | 0.000355 | 12.72 |
| DeepDoubleCvL | 0.004997 | 0.000478 | 10.45 | 0.000398 | 12.56 |

- Another observation: batch evaluation can bring substantial speed-up in some cases

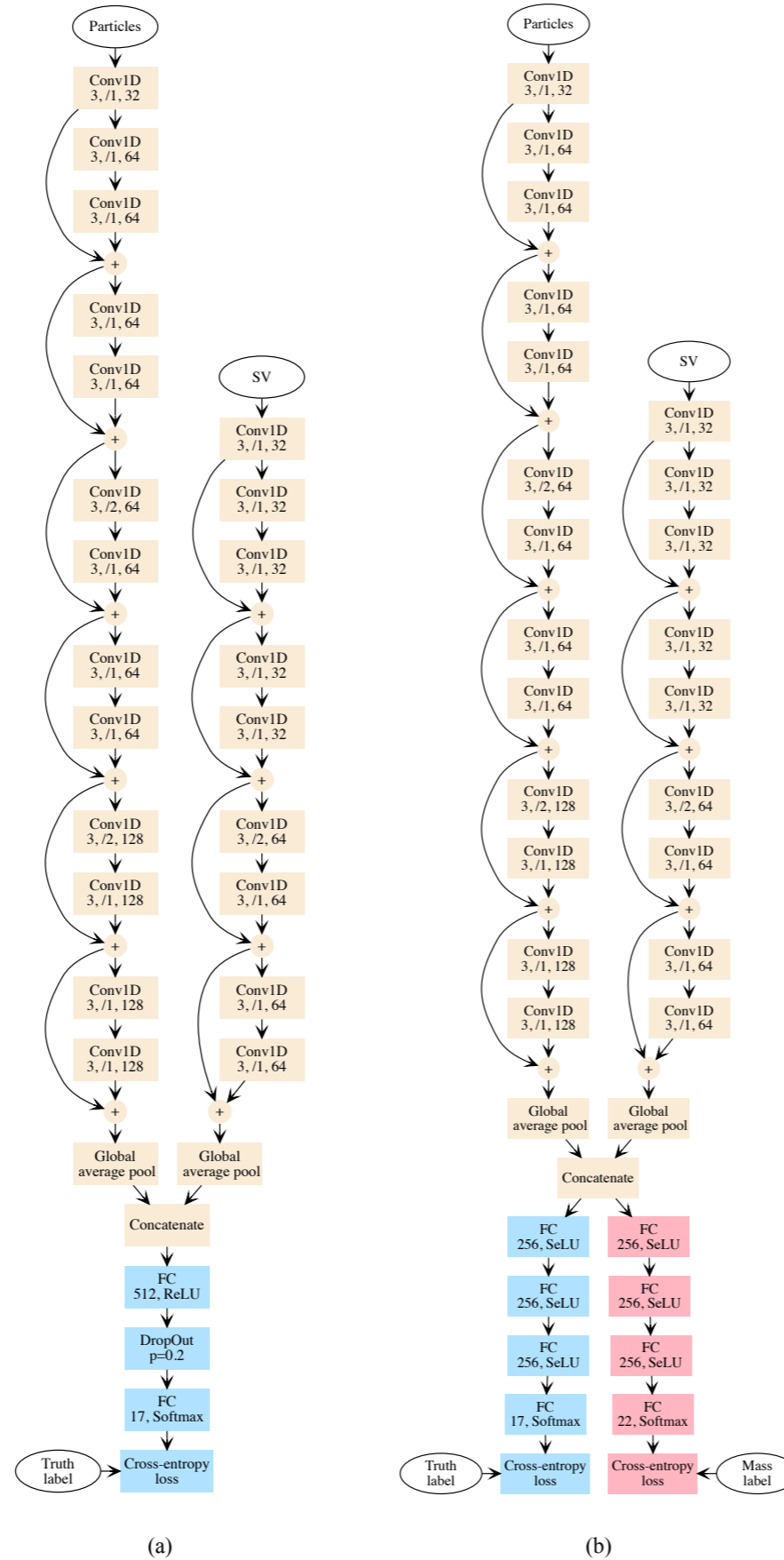  - right: DeepJet inference w/ and w/o batching

# SUMMARY

- A number of DL-based object reconstruction and identification algorithms have been developed in CMS over the past few years

  - significant improvement in performance

  - successfully applied to several challenging analyses and led to very competitive results

- The integration of DL frameworks into CMSSW is often a challenging task

  - multi-threading schemes

  - resource constraints (CPU time/memory)

- ML inference starts to become a sizable fraction of the event processing time

  - crucial to investigate how to accelerate ML inference

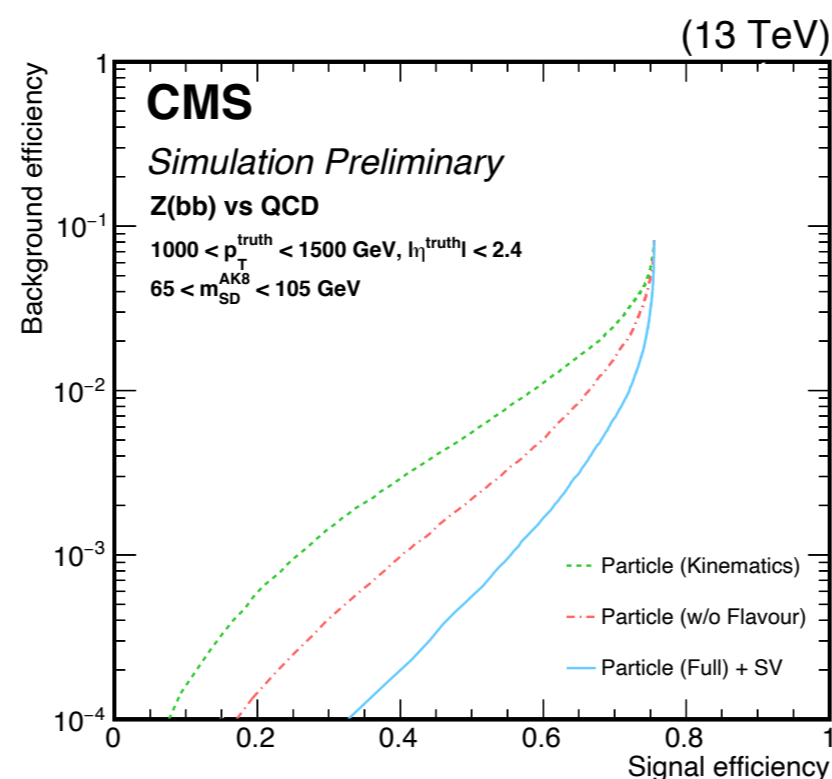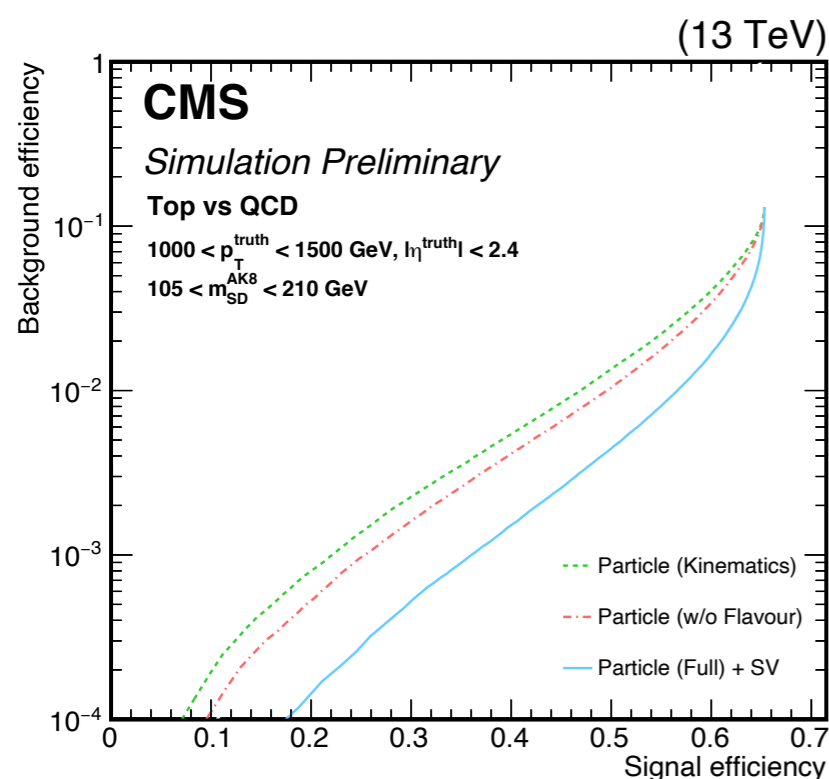    - e.g., new frameworks, new DNN architectures, new hardware, etc...

# BACKUPS

(a)  (b)

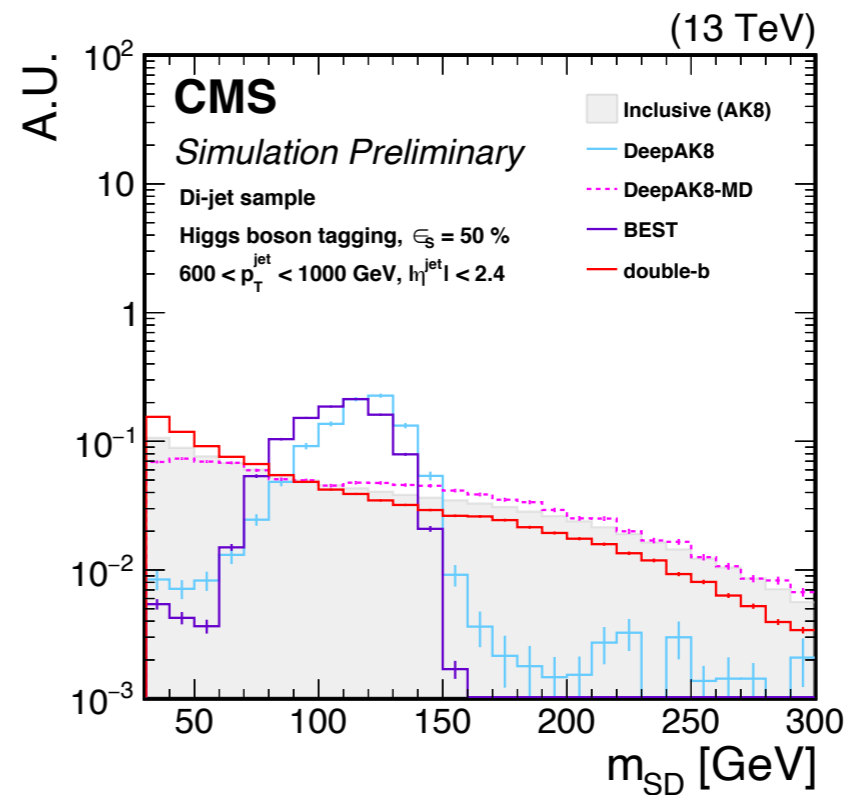*ML Inference in CMSSW - November 15, 2019 - Huilin Qu (UCSB)*

CMS-PAS-JME-18-002
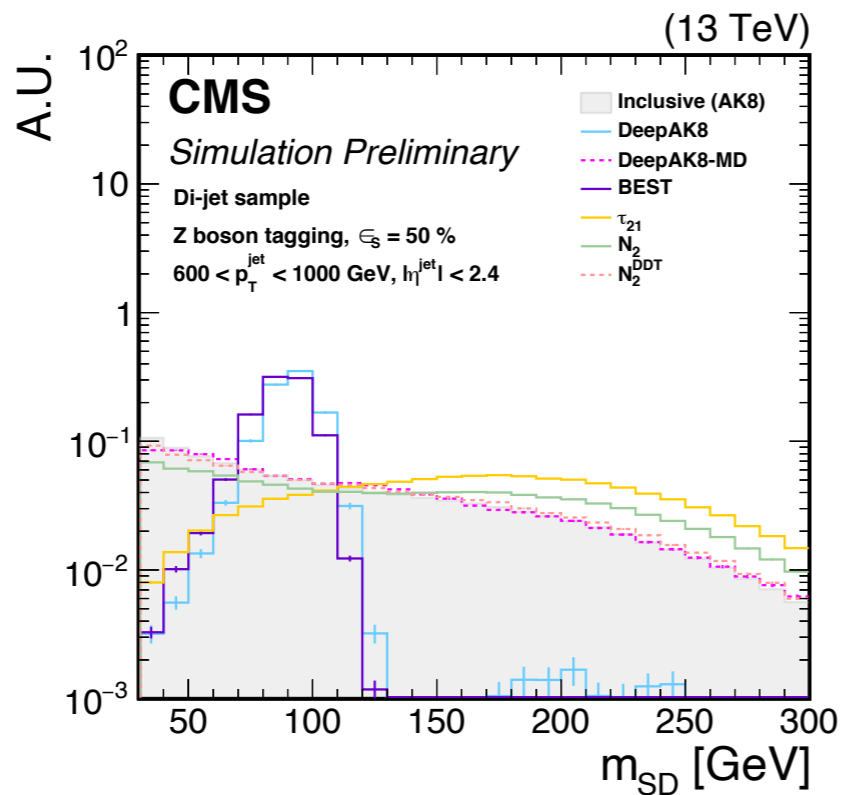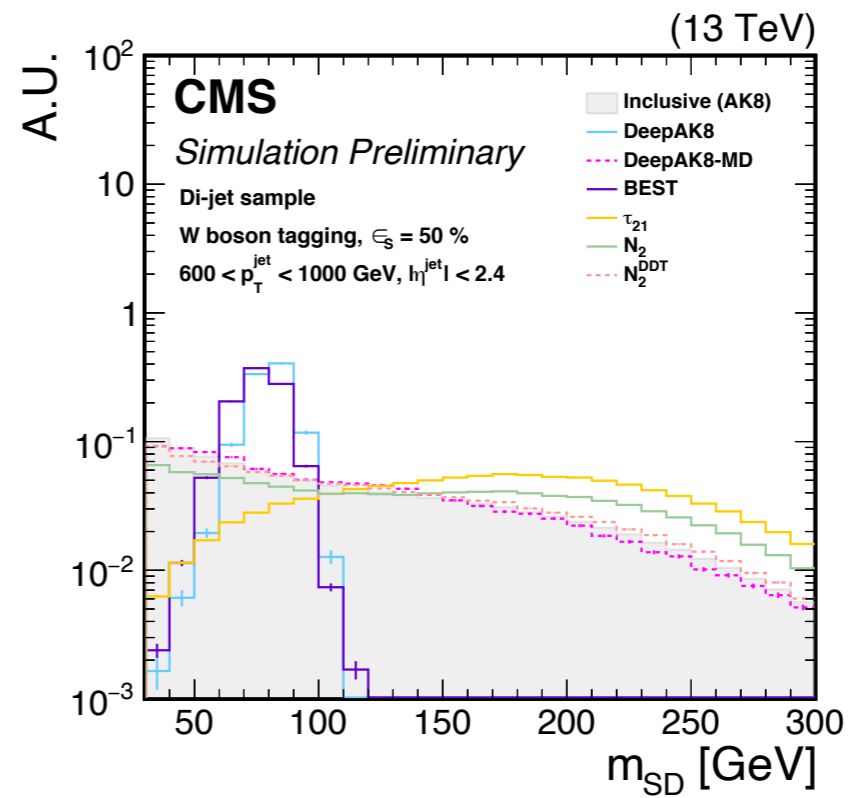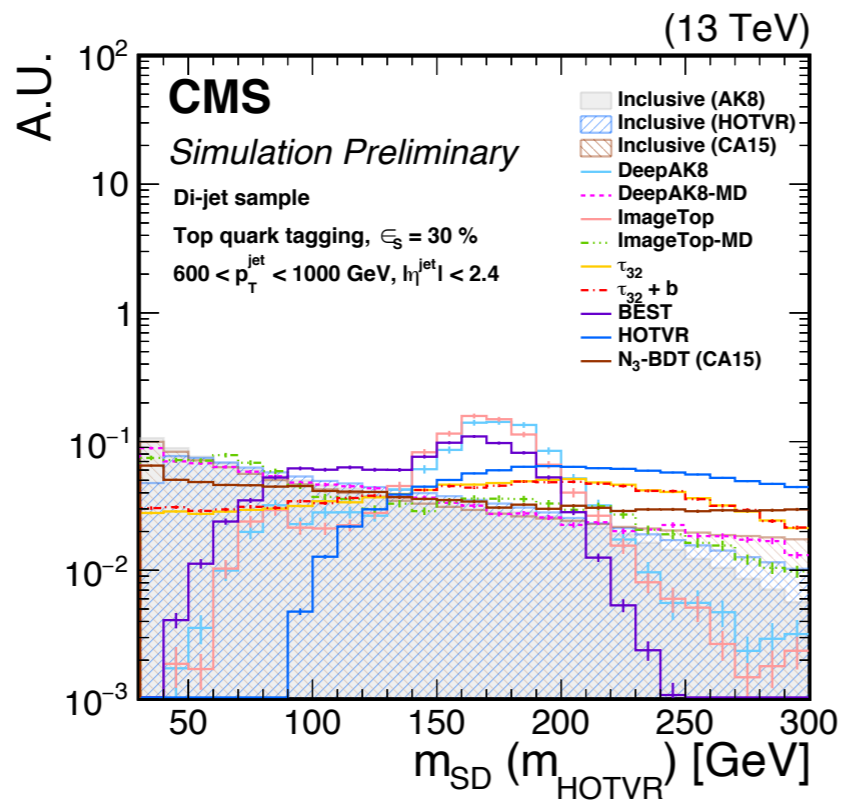
# ABLATION STUDY OF DEEPAK8

- DeepAK8 shows substantial gain compared to traditional approaches    *CMS-PAS-JME-18-002*

- To understand the main sources of the improvement, alternative versions of DeepAK8 were trained using a subset of the input features

  - Particle (kinematics): only kinematic info of PF candidates

    - four momenta, distances to the jet and subjet axes, etc.

  - Particle (w/o Flavour): adding experimental info

    - charge, particle identification, track quality, etc.

  - Particle Full + SV (the full DeepAK8): adding features related to heavy-flavour tagging

    - track displacement, track-vertex association, SV features, etc.

# MASS DECORRELATION

# DEEPTAU: ALGORITHM

- Input variables
  - **1 global event variable**: the average energy deposition density ($\rho$)
  - **42 high-level variables** that are used during tau reconstruction or proven to provide discriminating power by previous tau POG studies
  - For each candidate reconstructed within the tau signal or isolation cones, information about 4-momentum, track quality, relation with the primary vertex, calorimeter clusters, and muon stations is used, if available:
    - From **7 to 27 variables** (depending on the candidate type) for each **particle flow candidate**
    - **37 variables** for each fully reconstructed **electron candidate**
    - **37 variables** for each fully reconstructed **muon candidate**

- Candidates belonging to the inner and outer cones are separated and split into two grids with $\eta \times \varphi$ cell size of $0.02 \times 0.02$ ($0.05 \times 0.05$) for the inner (outer) cone

- Network architecture:
  - High level variables and each input cell are pre-processed by a few fully connected dense layers
  - For the inner (outer) grid, the pre-processed cell data are fed into 5 (10) 2D convolutional layers with $3 \times 3$ window size, which result in 64 features that are passed to the next step
  - All features from previous steps are combined and passed through 5 dense layers
  - Probabilities of the reconstructed tau candidate being electron, muon, quark or gluon jet, or hadronic tau are estimated by the 4 NN outputs

*CMS-DP-2019-033*

## Distribution of the visible $\mu\tau$ mass for 2018 data

Event selection:
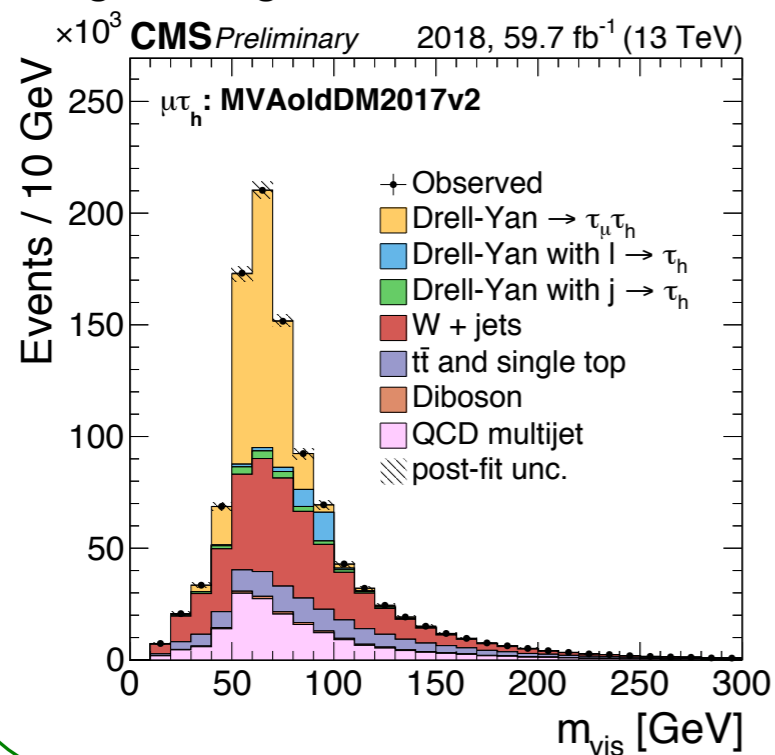- well identified and isolated muon with $p_T > 25$ GeV, $|\eta| < 2.4$, $|dz| < 0.2$ cm
- tau candidates with $p_T > 20$ GeV, $|\eta| < 2.3$, $|dz| < 0.2$ cm
- $\mu\tau$ pair with an opposite charge and $\Delta R(\mu, \tau) > 0.5$

- Contribution from all SM processes (except QCD) are modelled by MC simulation
- QCD estimated from a sideband region in data

**Selection using discriminators from JINST 13 (2018) P10005:**
- Tight WP against jets
- VLoose WP against electrons
- Tight WP against muons



With DeepTau selection, the yield from **genuine $\tau_h$** increases by ≈ 20%, while yield from **fakes** decreases by ≈ 23%

*In both plots modelled contributions are fit to the data*

**Selection using DeepTau IDs:**
- Tight WP against jets
- VVLoose WP against electrons
- VLoose against muons

CMS-DP-2019-033