# Update on Spark Workflow @ Vanderbilt

Andrew Melo
CMS Big Data Meeting
Mar 21, 2018

# Spark @ Vanderbilt

- Vanderbilt has two "Big Data" clusters - one currently dedicated to CMS

  - Mostly purchased, some hardware "harvested" from EOL compute racks

- Slave machines run Mesos

  - Both Spark and Jupyterhub allocate from Mesos

  - Would like to backfill idle resources with regular CMS jobs

- 160/50TB raw/usable HDFS storage, adding another 144/48TB

- Primary user interface is via Jupyterhub, CLI/programmatic access also available
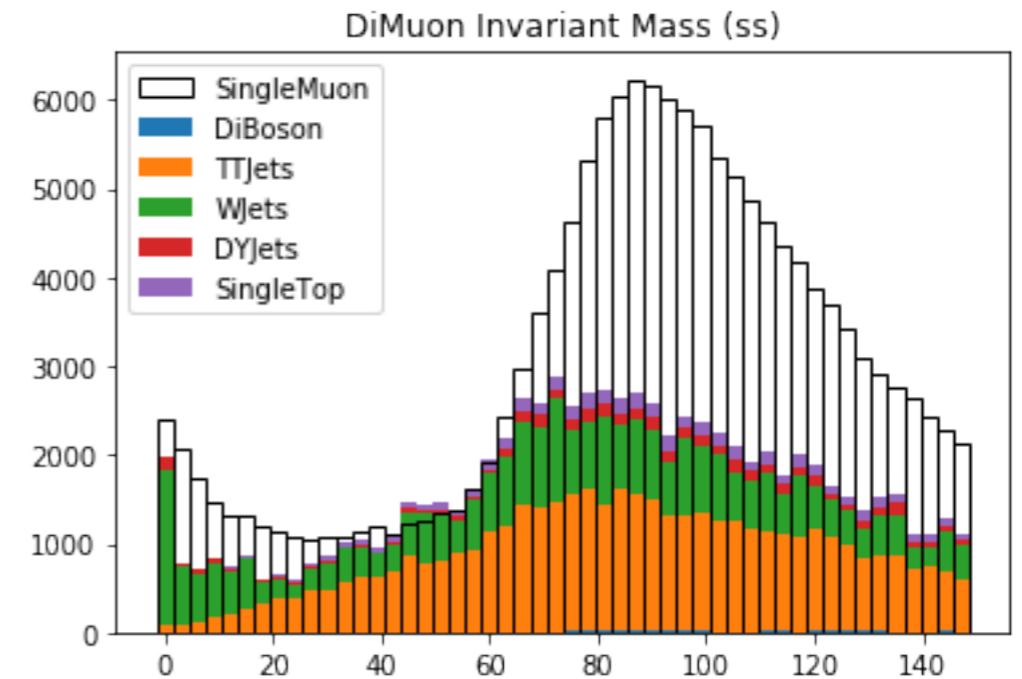
# Goals

1. Reproduce AN-17-142 with Spark+Jupyter

2. Demonstrate & train others on this technique
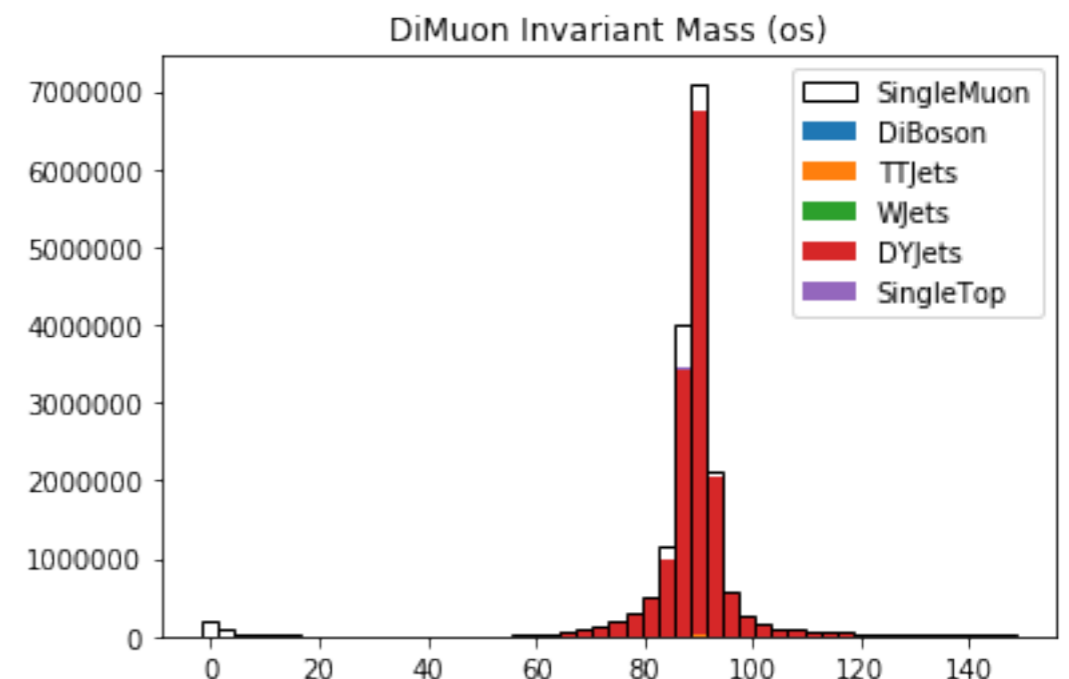
3. Get ≥1 "not savvy" user

# Progress

- All the parts are there! Now it's just polish & quality-of-life improvements

- Demoing one control region of AN-17-142 to my physics group today*

  - Will upload this afternoon

- Want to update last fall's HATS**

**\* Should've been last week, master hosts got powered down by mistake**

**\*\* https://github.com/FNALLPC/spark-hats**



**Produced from 931M events in ~90 secs**

# Next Steps/Questions

- Finish converting AN-17-142 2x

  - Short demonstration of one CR ✓

  - Full reproduction of plots/tables ❌

- Install XRootD cache

- Streamline EOSConnector deployment

- Find a ~~victim~~ user

- Write a library to replace the (substantial) boilerplate?

- How do we make notebooks interoperable between sites?

  - Storage locations, etc

- How should we best advertise these resources to users?

- What are best practices for Mesos administration?

  - Backfilling idle resources, fair resource allocation between users

- Has PubComm updated their regulations?