

ML based Anomaly Detection at the LHC

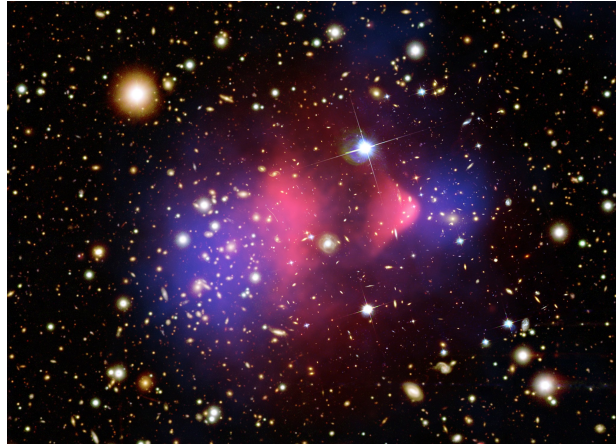
Oz Amram



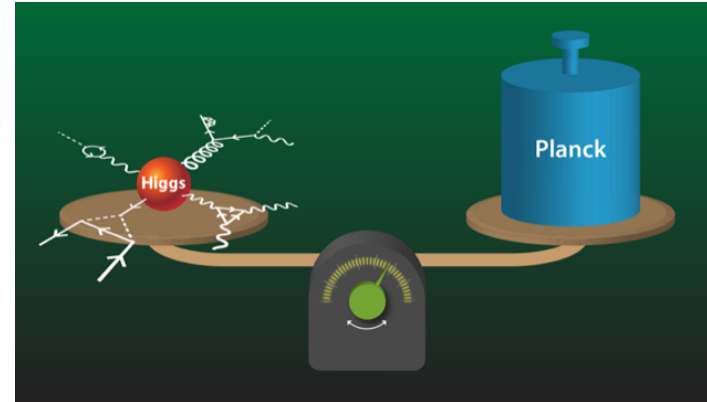
**U Chicago
Rising Stars Symposium
September 2021**

Lots of Questions

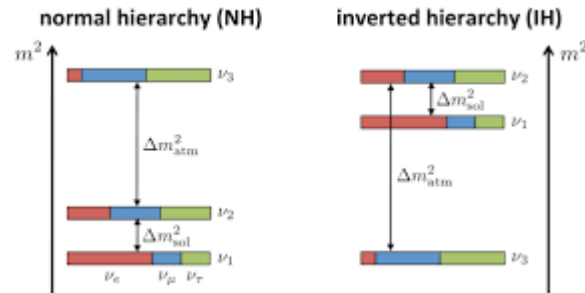
What is Dark Matter?



Why is the Higgs so light?



What is the origin of Neutrino Mass?



Baryogenesis?

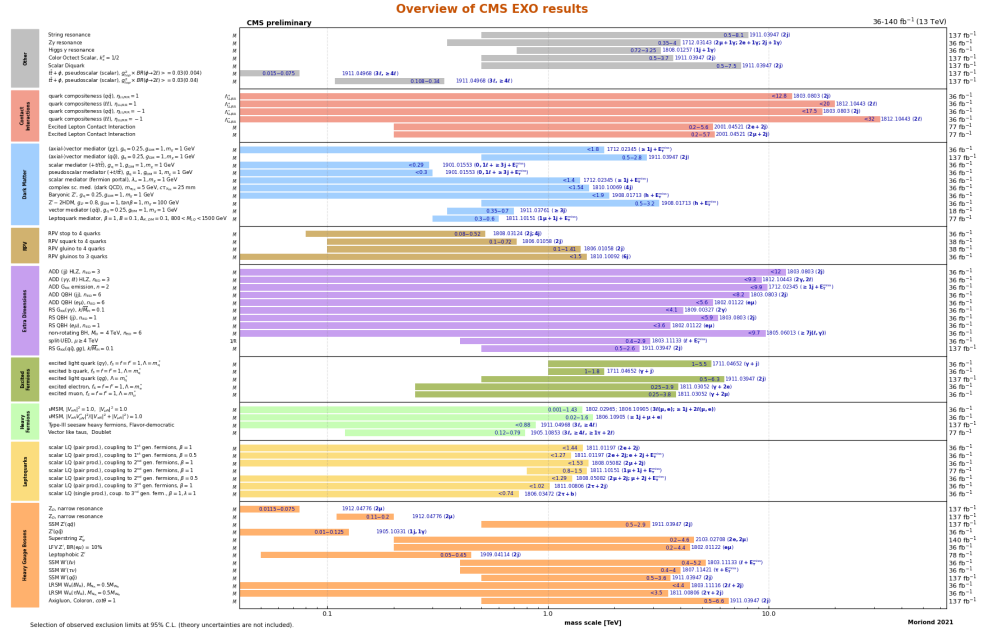
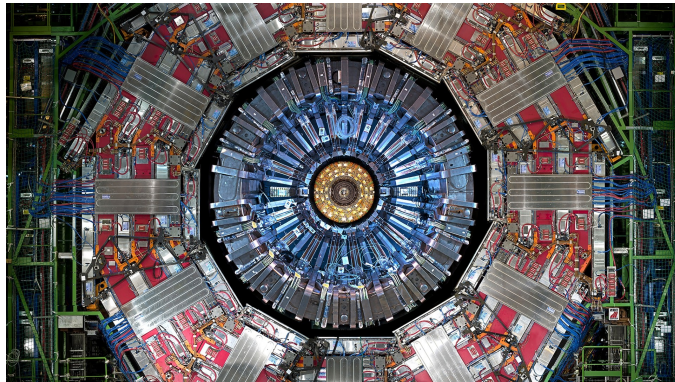
Grand Unification?

And many more...

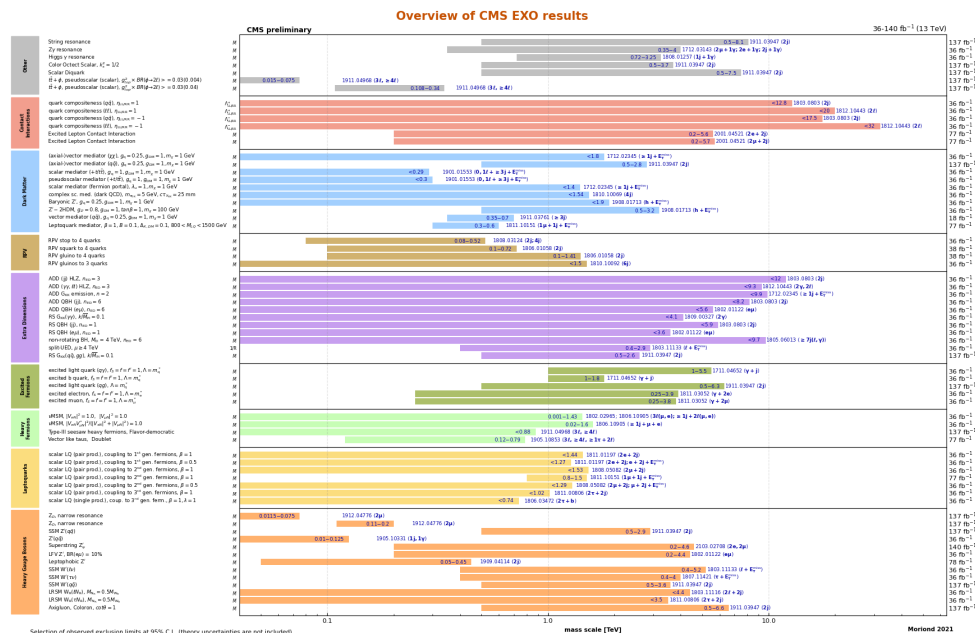
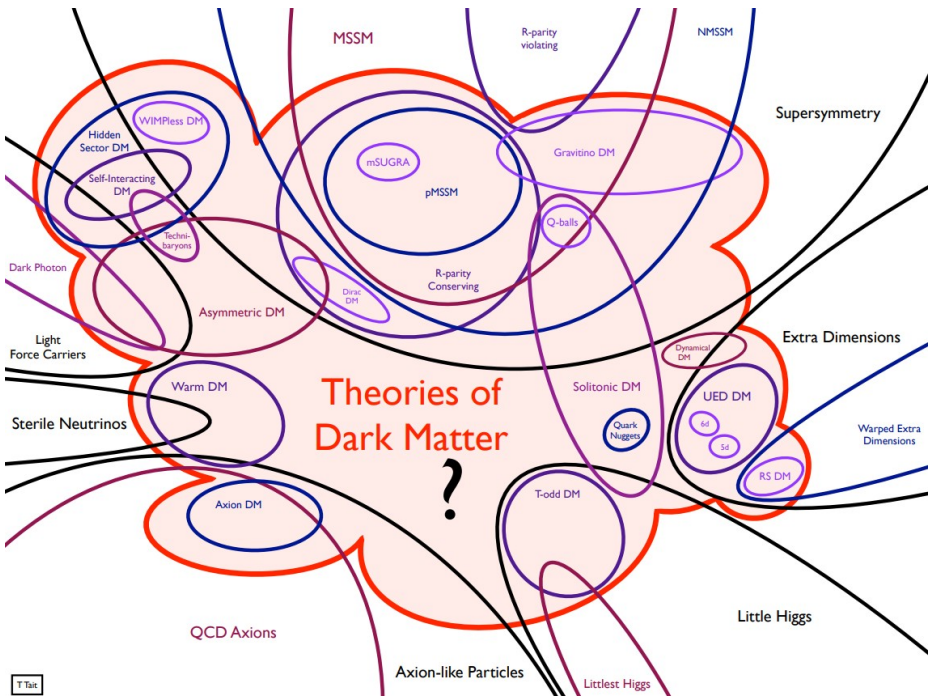
g-2?

Flavor Anomalies?

No clear answers from the LHC yet



But..



What if we aren't looking in the right places ?!



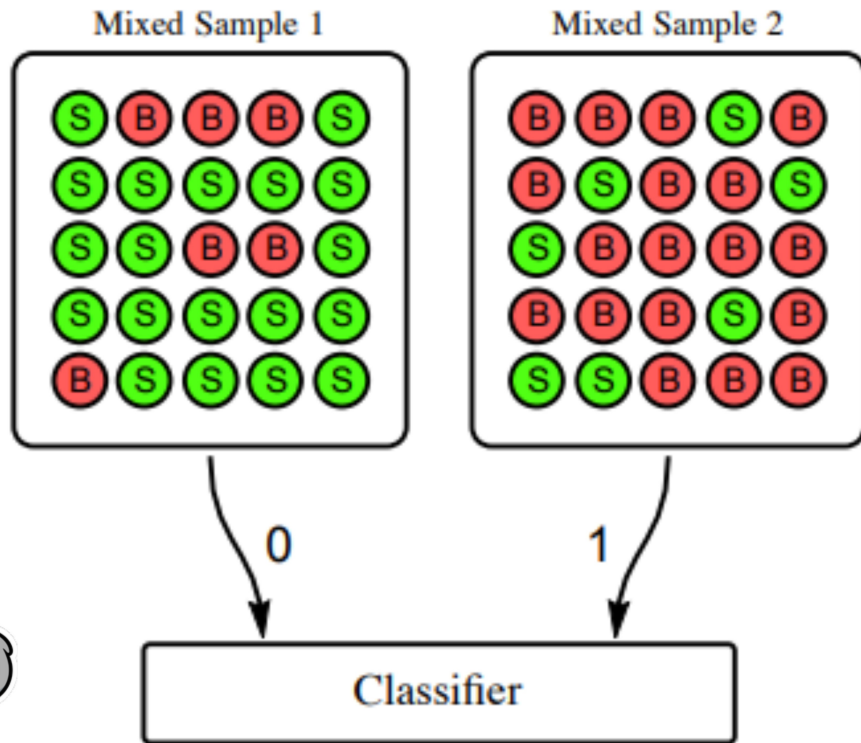
The Challenge

- How can we design searches with minimal assumptions but still have powerful sensitivity?
- New ideas in ML are enabling totally new search strategies!

Key Idea: Train directly on data!

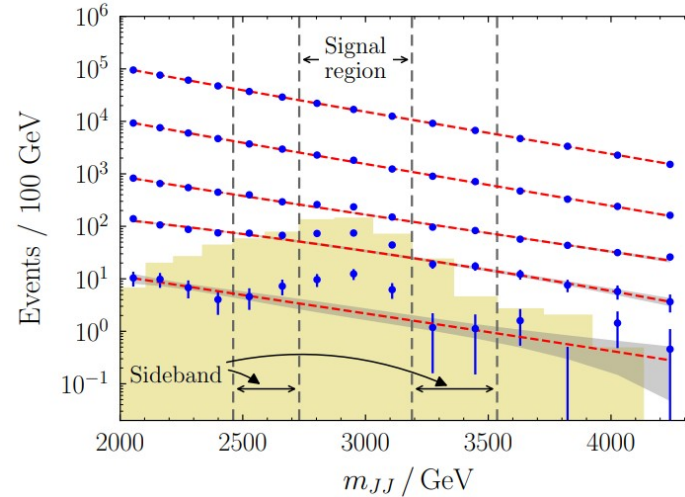
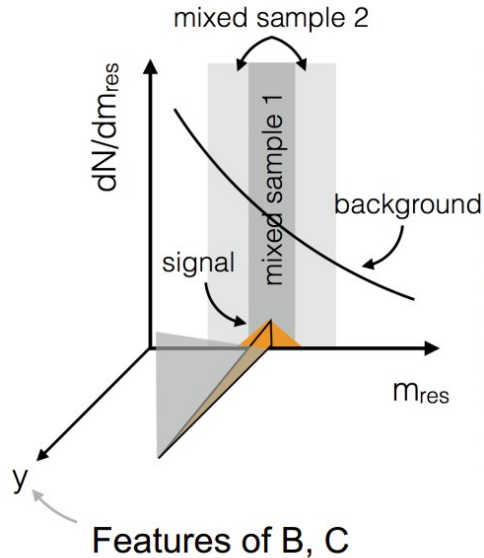
Classification Without Labels

Metodiev, Nachman & Thaler 1708.02949



CWoLa Hunting

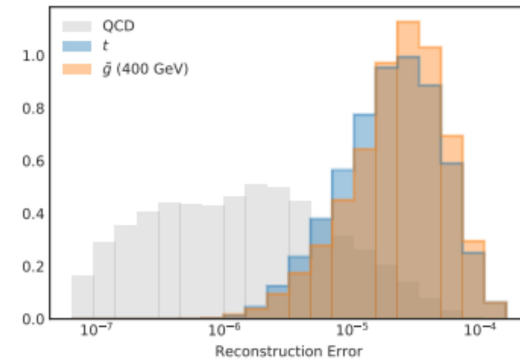
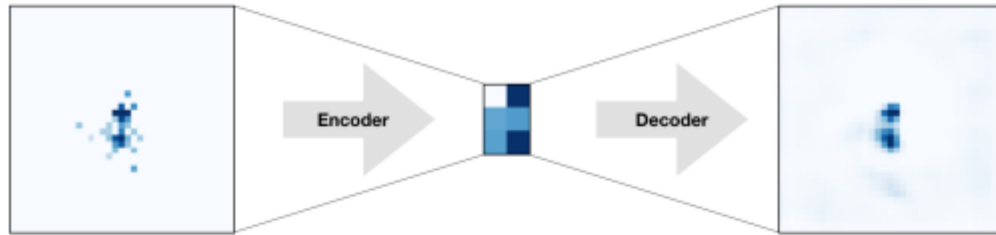
Colins, Howe & Nachman [1902.02634](#)



- Signal region = dijet mass window
- Train a classifier on signal region vs. others
- Select events & bump hunt



Anomaly Detection : Autoencoders



- Train a network to compress and decompress the data
- Can train directly on data, no labels needed
- Anomalous events should have a higher reconstruction loss

Drawbacks

- CWoLa Hunting
 - Worry about sculpting QCD dijet mass distribution
 - Apply to non-resonant signals?
- Autoencoders
 - Only 'learns' what QCD looks like
 - Room for improvement as a Sig vs. Bkg classifier

Tag N' Train (TNT)

- A method of training improved classifiers on data
- Assumptions:
 - Signal has **2 interesting objects** in it
 - One has a **starting classifier** for each object
 - Signal-like features in background events are uncorrelated between the 2 objects

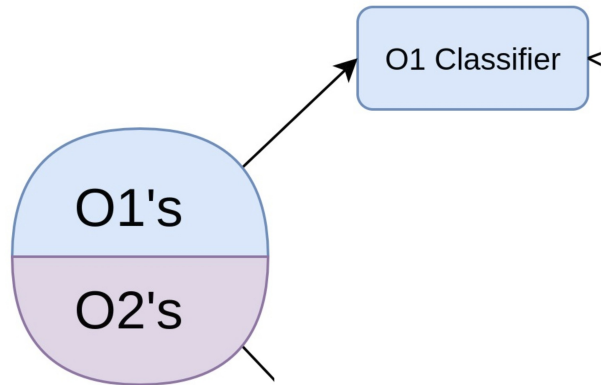


Tag with a weak classifier **N' Train** a better one!

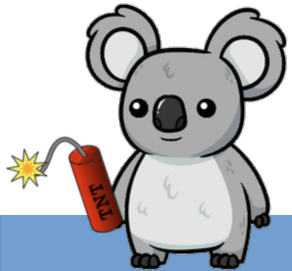
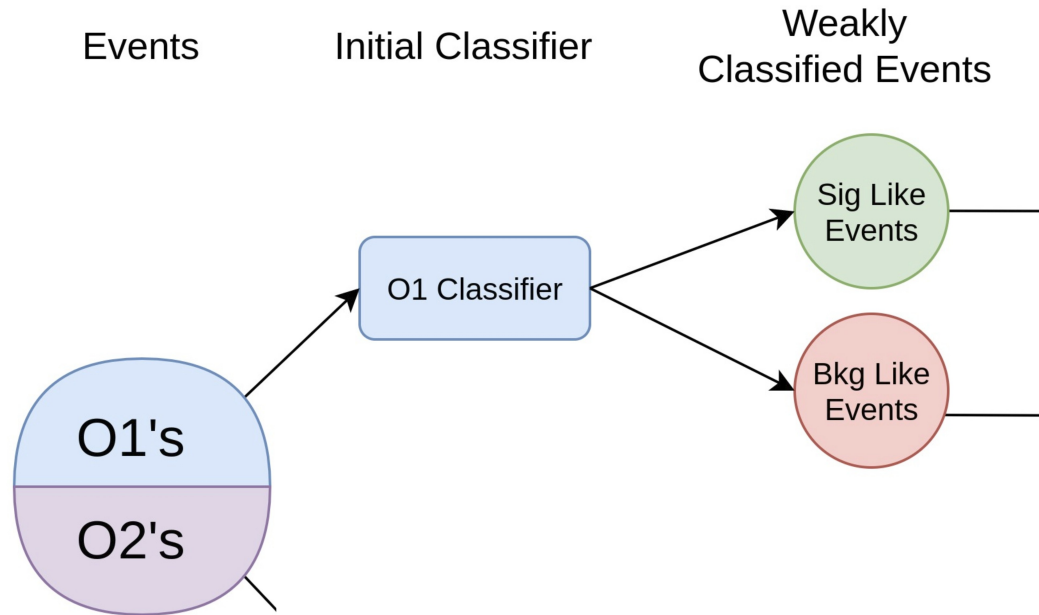
Tag N' Train

Events

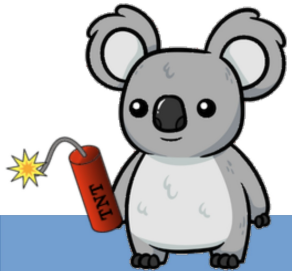
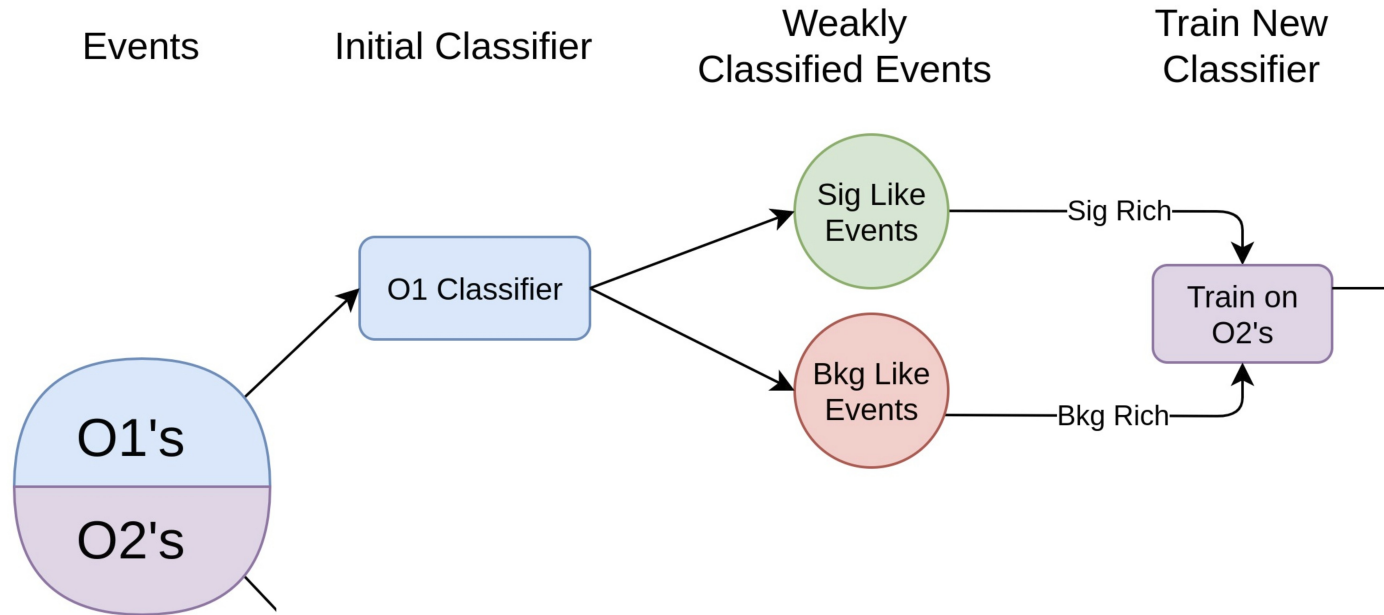
Initial Classifier



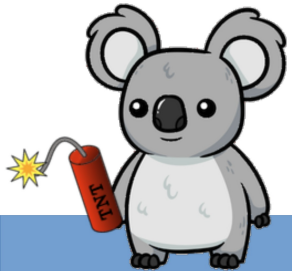
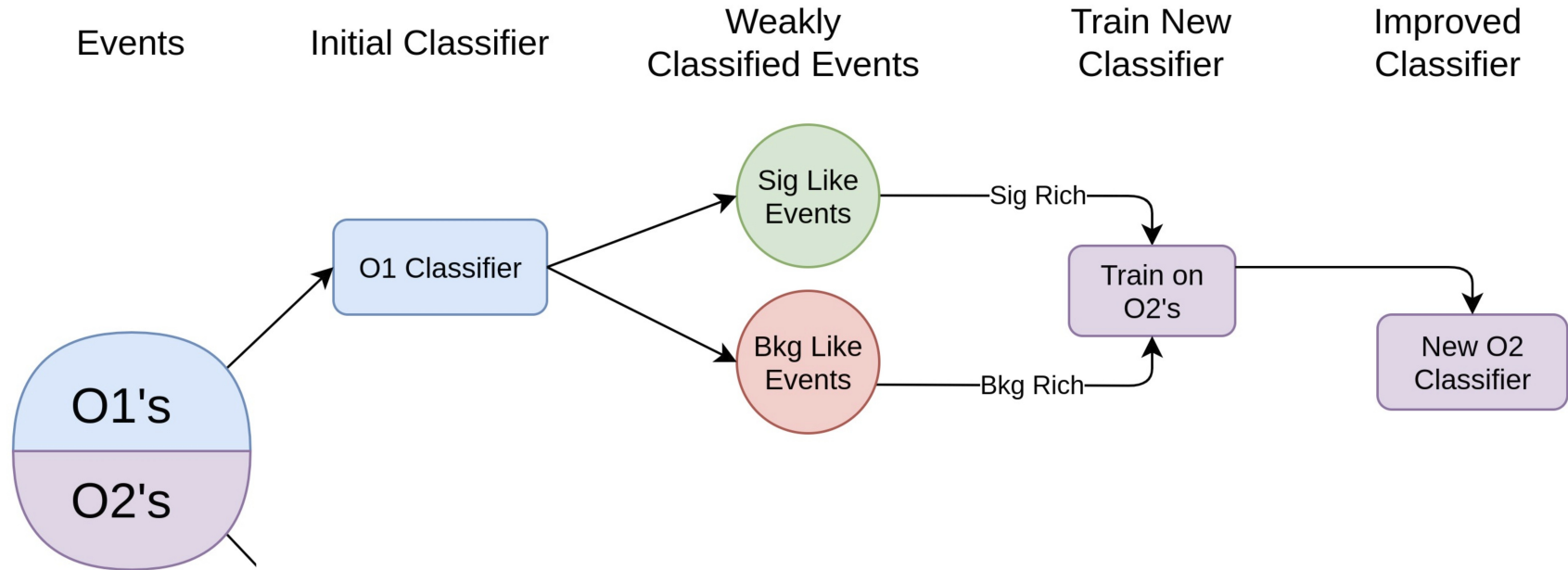
Tag N' Train



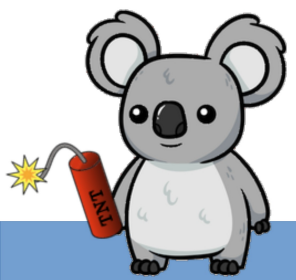
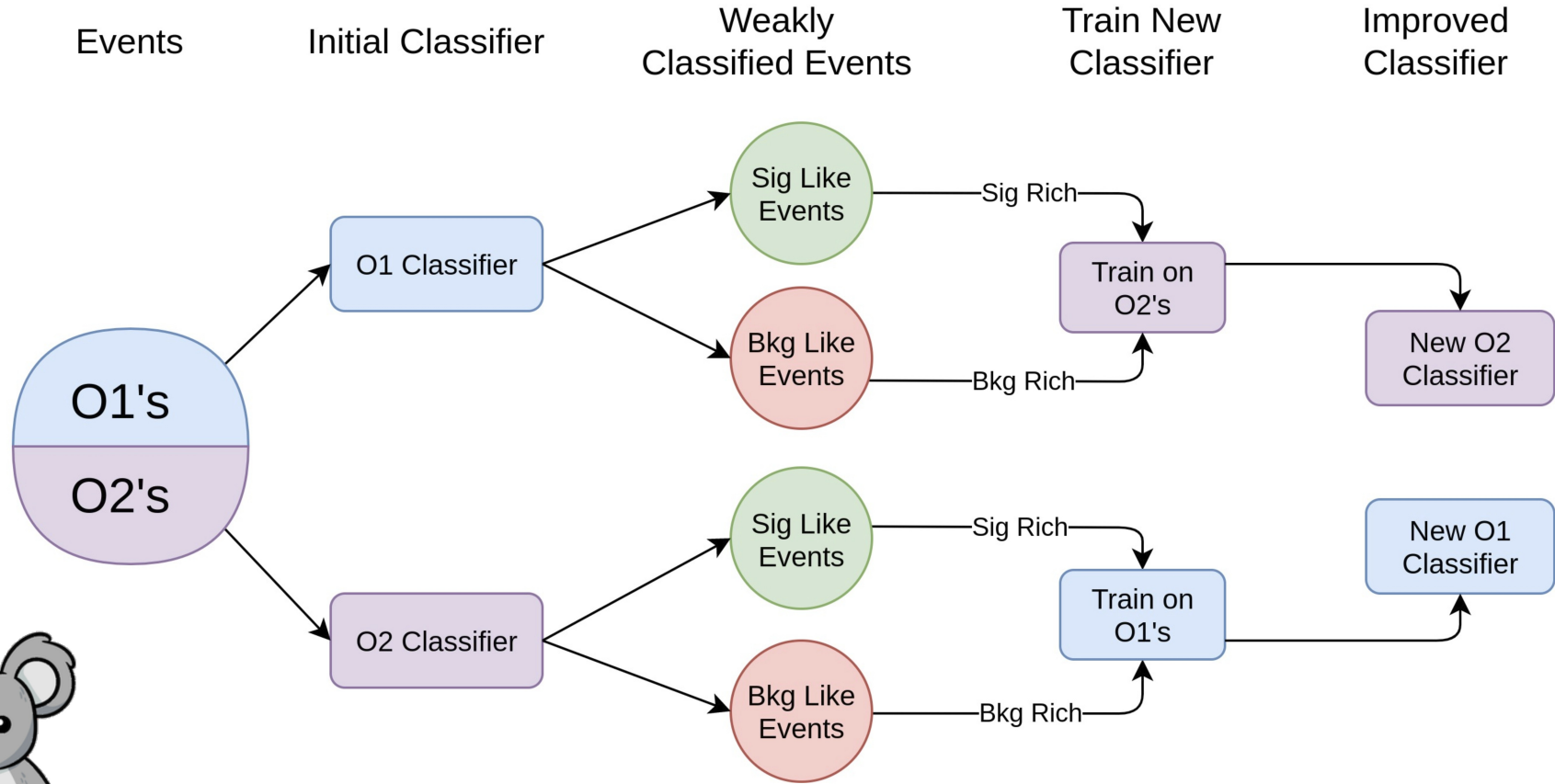
Tag N' Train



Tag N' Train

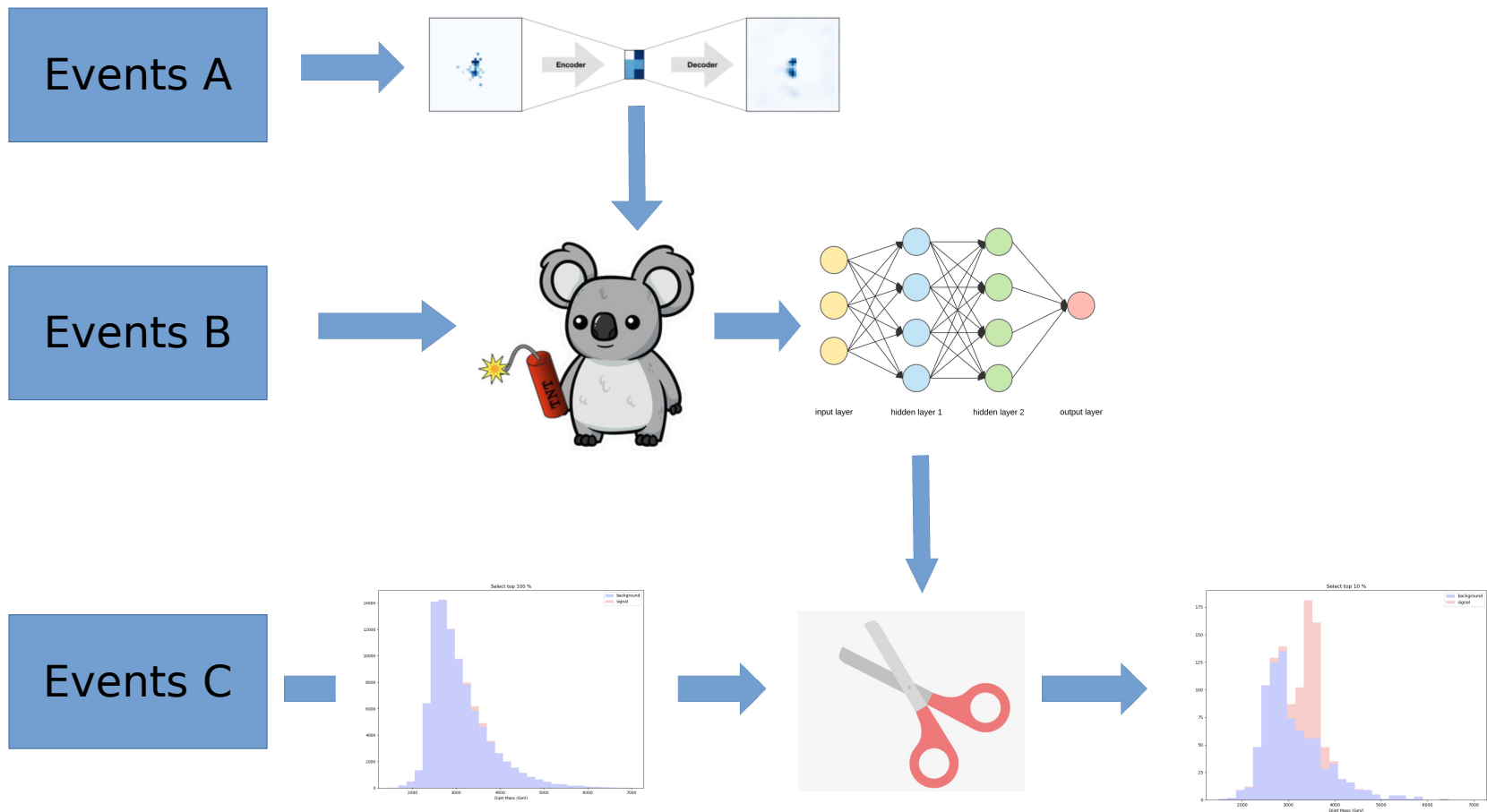


Tag N' Train

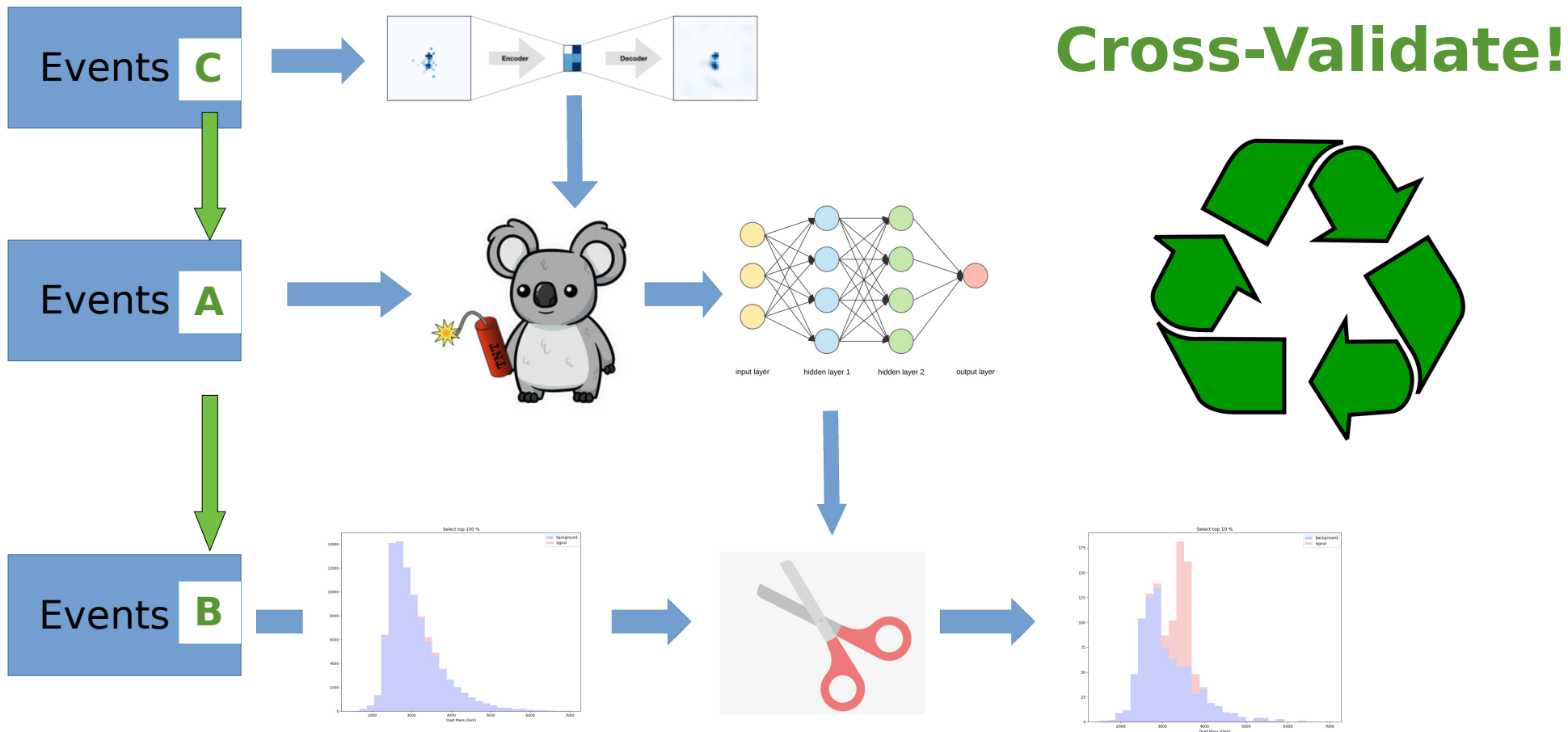


Dijet Anomaly Search

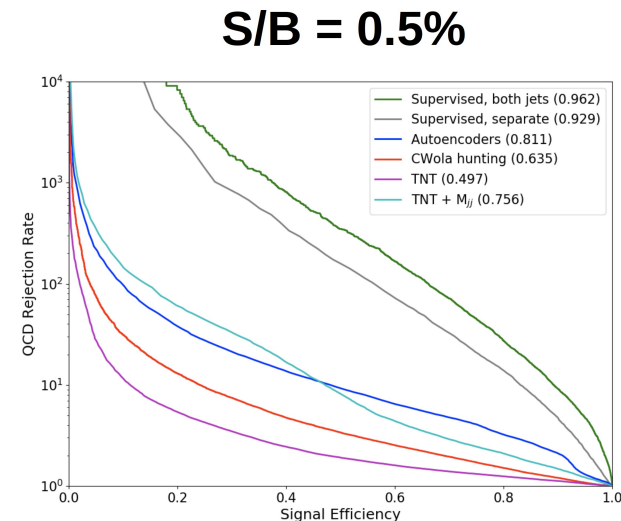
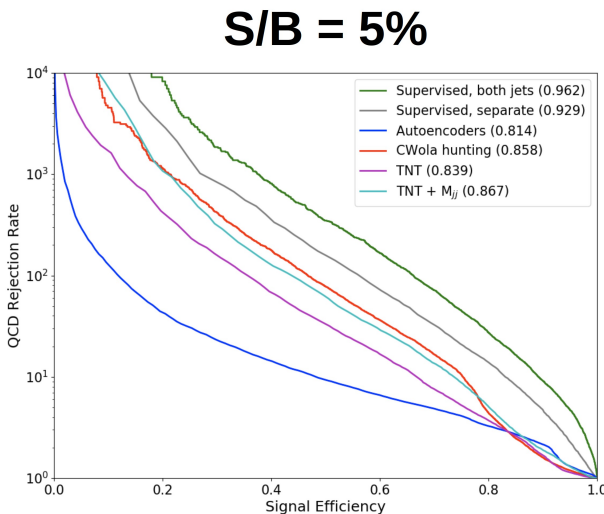
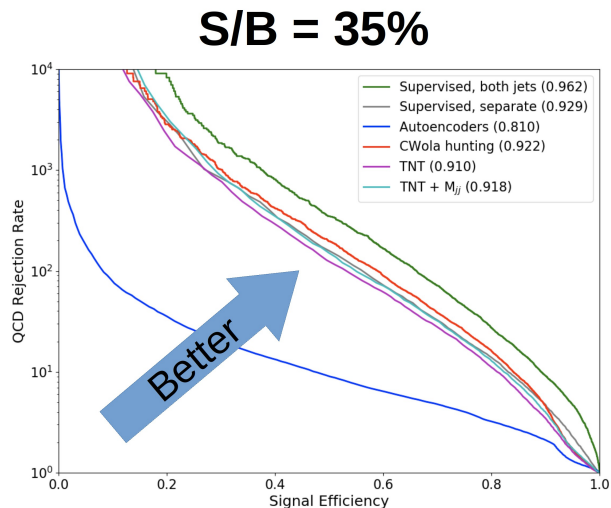
Applying TNT to a Resonance Search



Applying TNT to a Resonance Search



Classification Performance

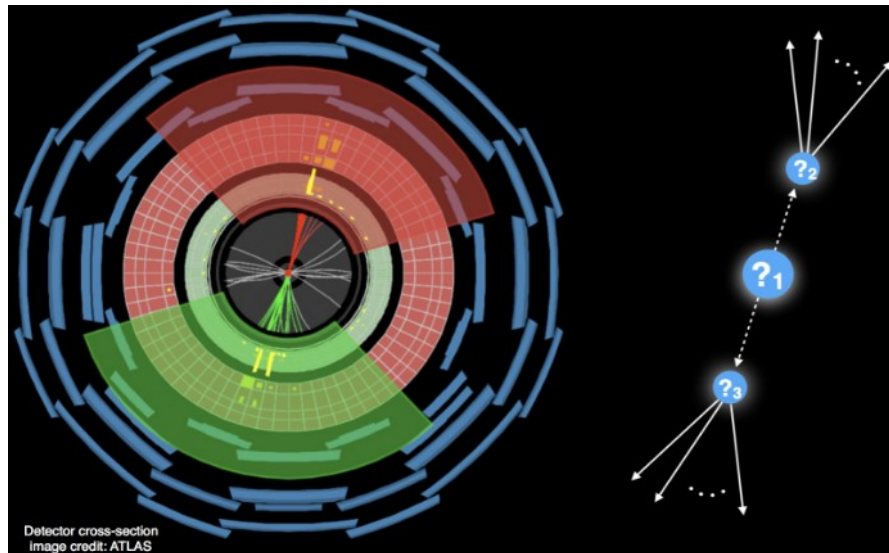


Less Signal

- CWoLa based methods approach **supervised** case when lots of signal
- **Autoencoders** performance independent of signal
- **TNT** matches **CWoLa hunting** high/medium signal, better at low signal

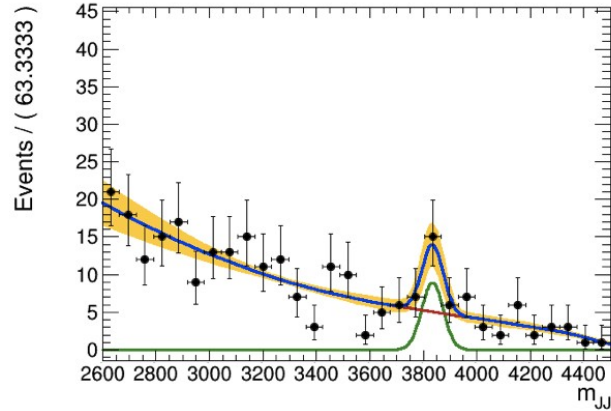
In Action: LHC Olympics 2020

- A competition to test out these new anomaly detection methods
- Blackboxes with:
 - 1M events, $R=1$ jet $p_t > 1.2$ TeV trigger
 - 4 vectors of all reconstructed particles
 - Mostly background + some hidden new physics (?)



[arXiv:2101.08320](https://arxiv.org/abs/2101.08320)

LHC Olympics Results

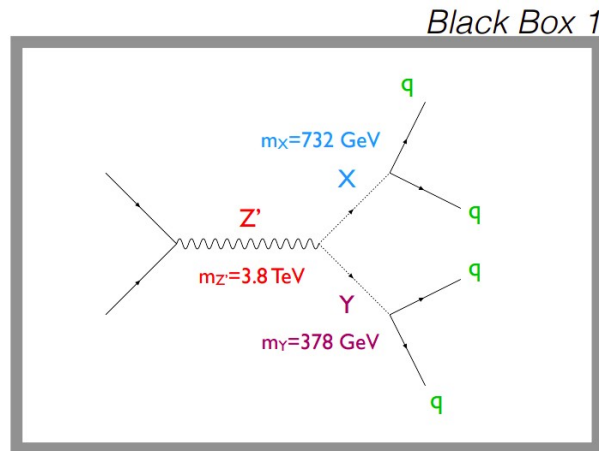
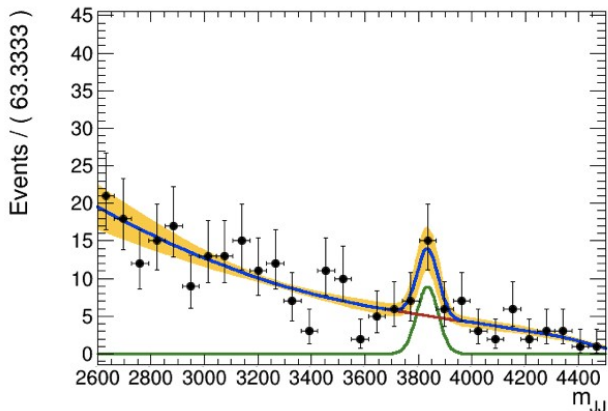


Black Box 1

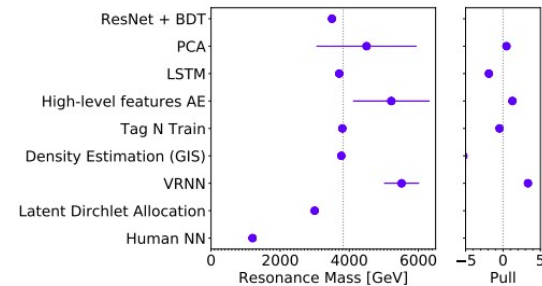


- TNT found a resonance at ~ 3800 GeV with 4σ evidence

LHC Olympics Results

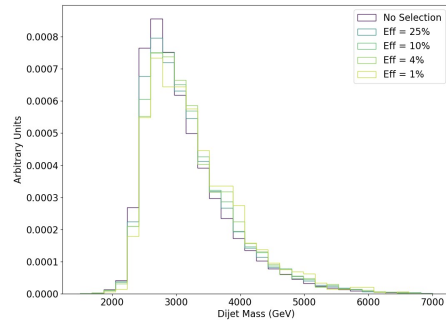


- TNT found a resonance at ~ 3800 GeV with 4σ evidence
- One of the few groups able to find the signal!

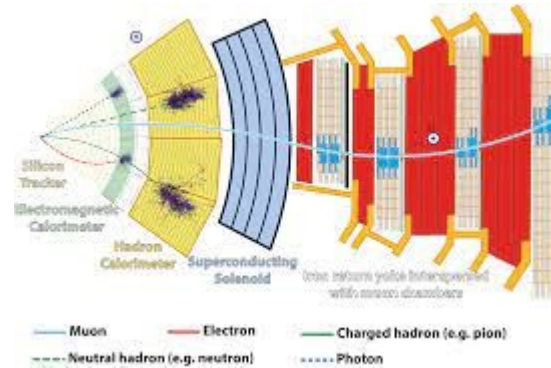


Many Challenges in Applying to Real Data!

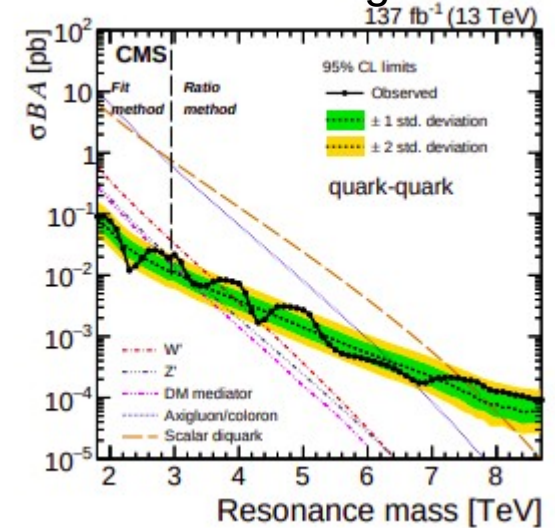
Ensure no mass sculpting



Don't "discover" a detector glitch!



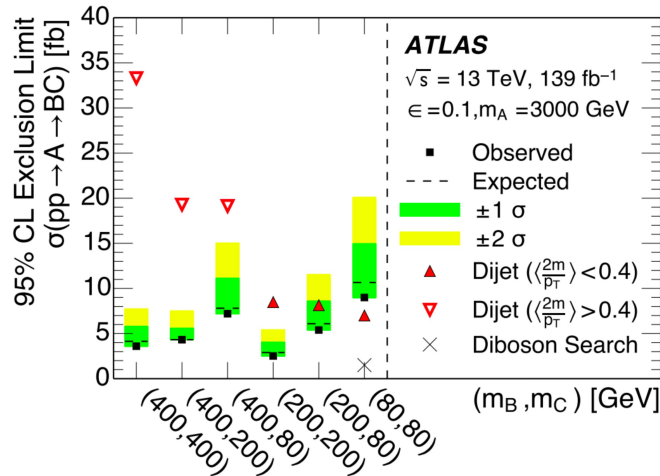
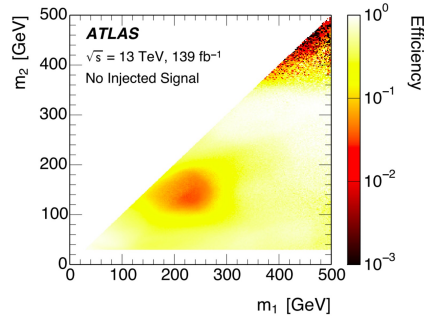
Limit setting ?!



Results on Data

ATLAS: 2005.02983

CMS



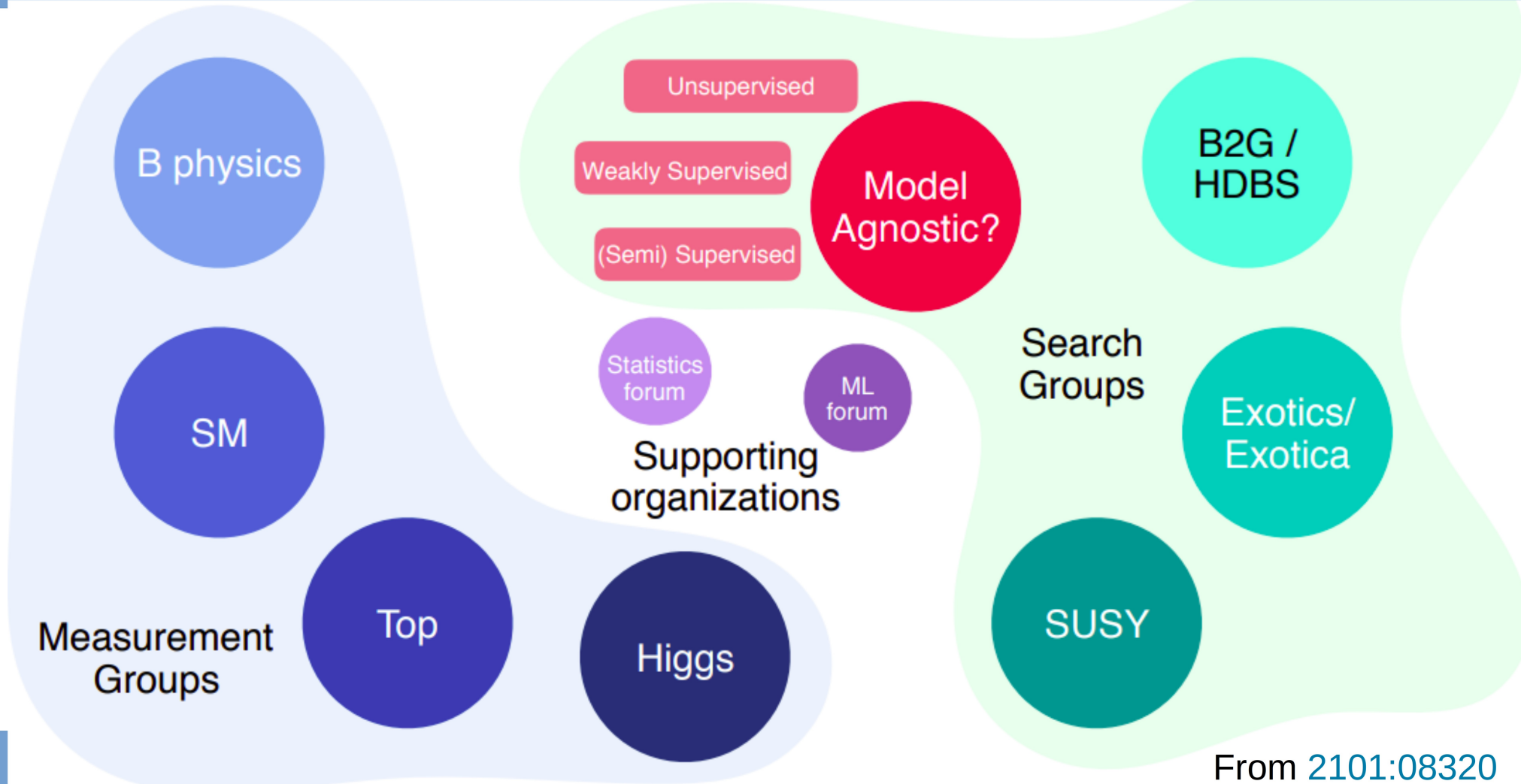
Just the Beginning...

- New techniques!
 - New ideas innovations from the ML side
 - Hybrid approaches with traditional searches
- New Searches!
 - Other anomalies besides jets with substructure
 - Non-resonant searches
- Do it fast!
 - Incorporate these ideas into triggers
 - Recently announced [Anomaly Detection at 40 MHz](#) challenge!

Current LHC Analysis Group Organization

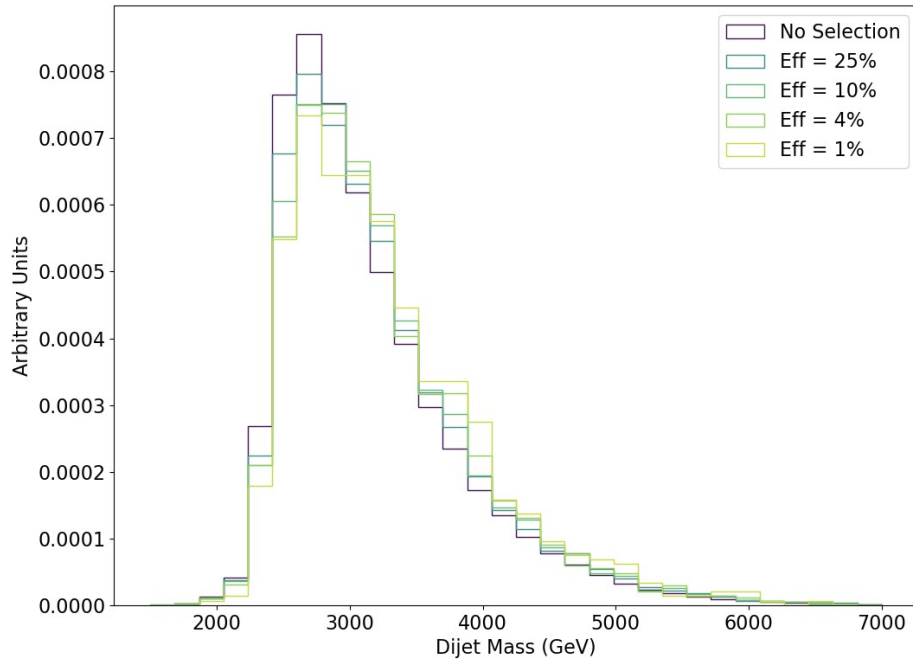


In 10 Years?



Backup

Dijet Mass Sculpting



- No sculpting of dijet mass!
- Decorrelation methods also possible with TNT
 - p_T reweighting tried, found no difference

TNT Technical Details

- 2 objects: heavy jet and light jet in event
- TNT Classifiers and autoencoders are CNN's based on jet images
- Top 20% 'sig-like', bottom 40% 'background-like'
 - Optional: require signal events in dijet mass window
- Combine 2 classifiers into 1
 - Require both jet's scores be in top X% of scores

Trade Offs*

(V)AE's



- + Performance indep. of amount of signal
- + Minimal assumptions
- Inherently ‘anti-QCD’ rather than a ‘pro-signal’

CWoLa Hunting



- + Great performance for large to medium signals
- + Can do full-event classification
- Assumption: resonant signal
- Must fully decorrelate features with M_{jj}

TNT

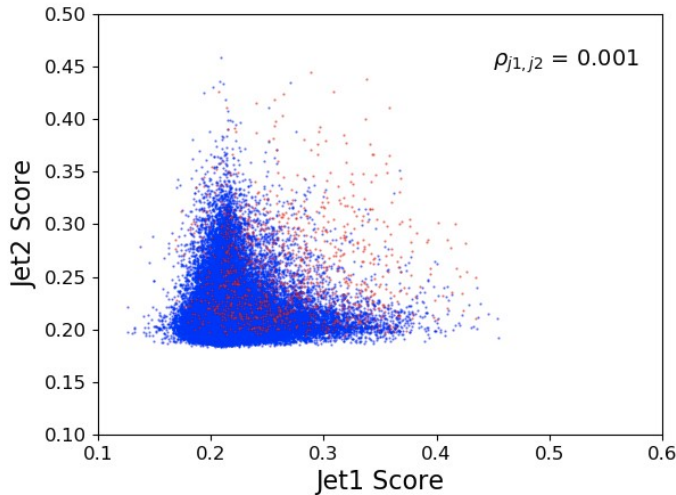


- + Great performance for medium/large signals and maintains performance for smaller signals
- + Mass sculpting mitigation possible
- Requires a starting classifier
- Assumption: Signal has 2 interesting objects

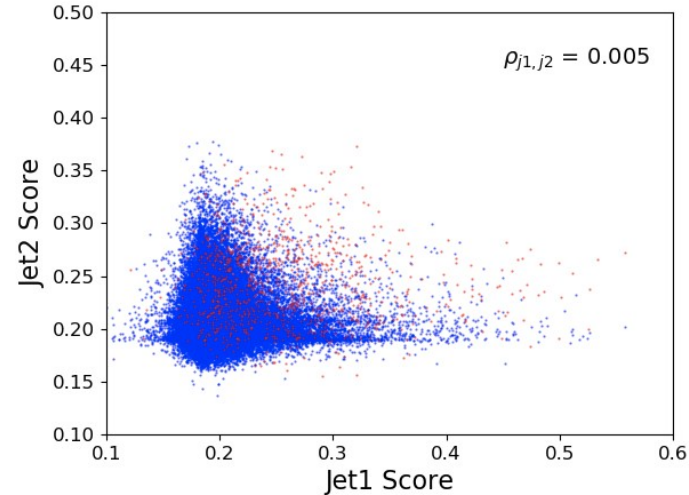
*Of course there are other interesting techniques with different trade offs too

Assumption: Correlations

“Pure” CWoLa



Tag N' Train



- Key assumption: Anomalous features of background events are uncorrelated
- Empirically (?) seems to hold