

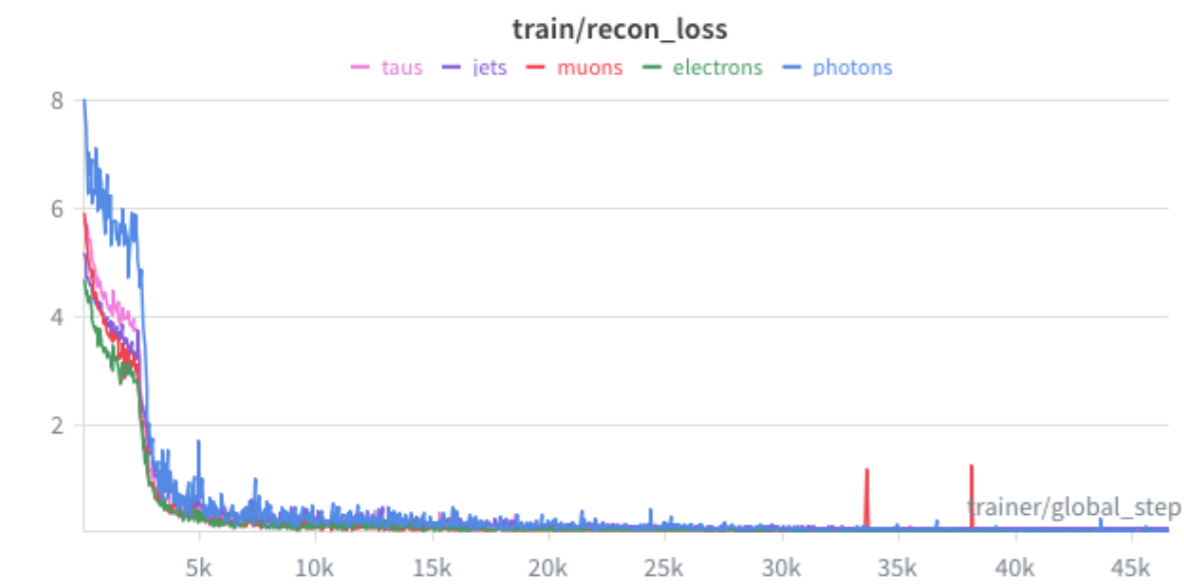
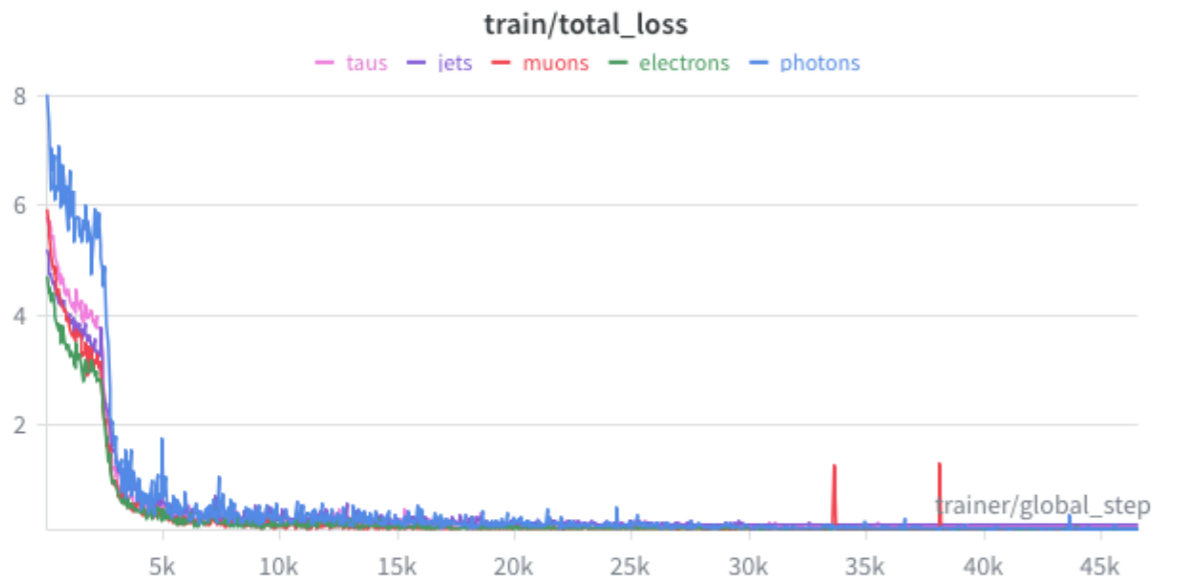
Treasure pipeline discussion

Merve Nazlim Agaras

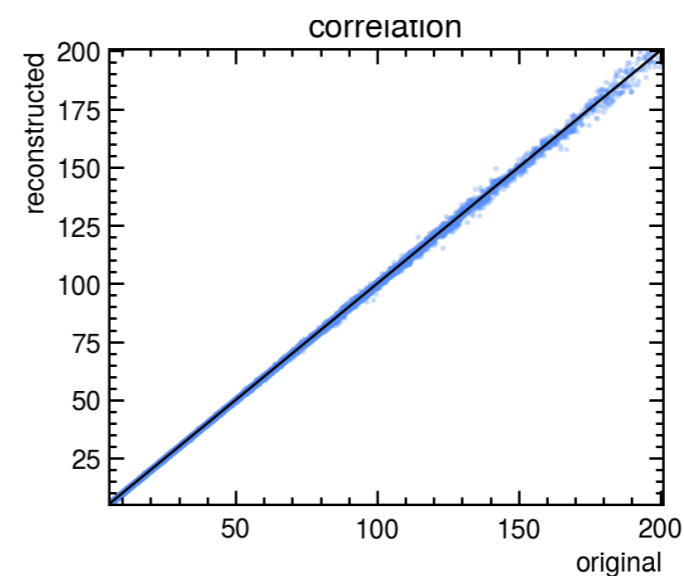
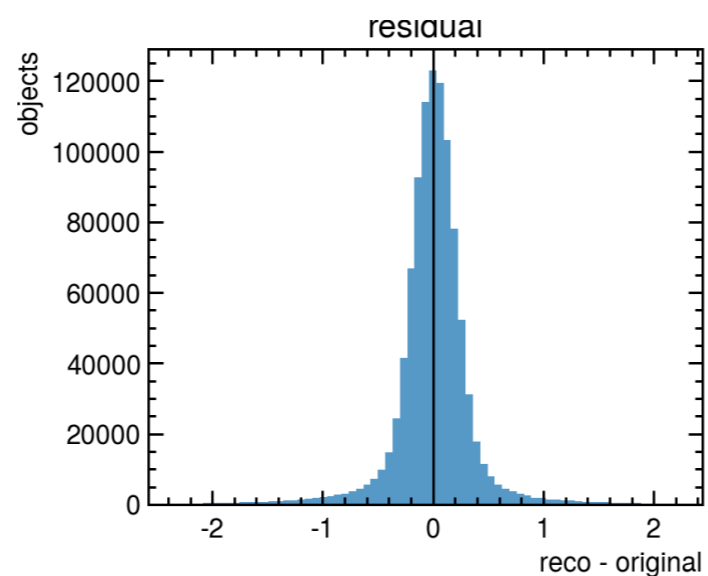
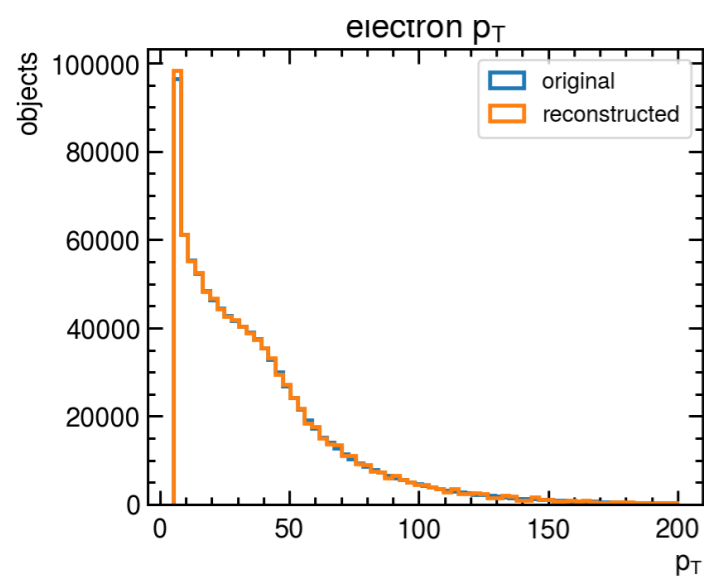
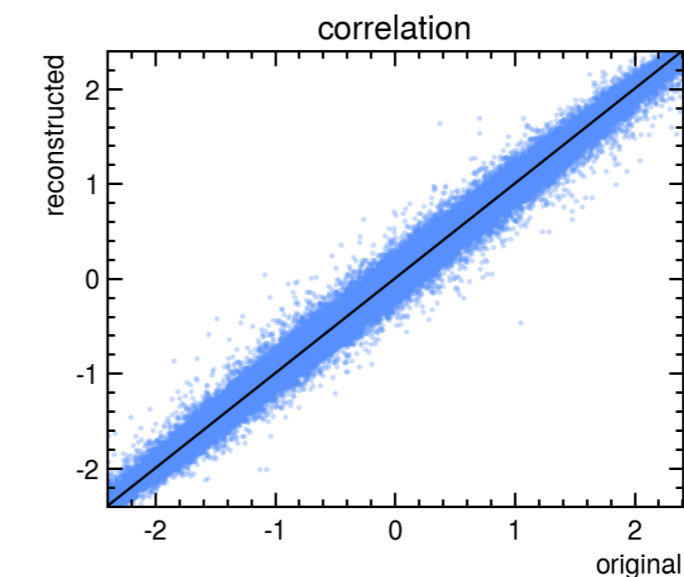
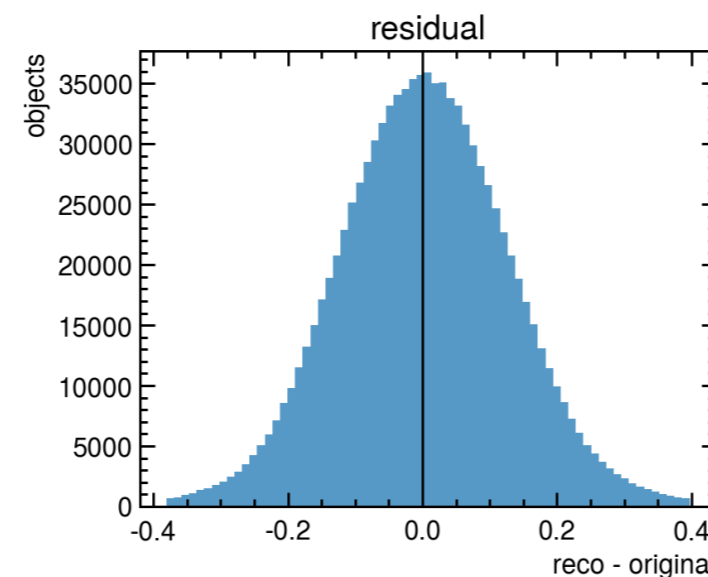
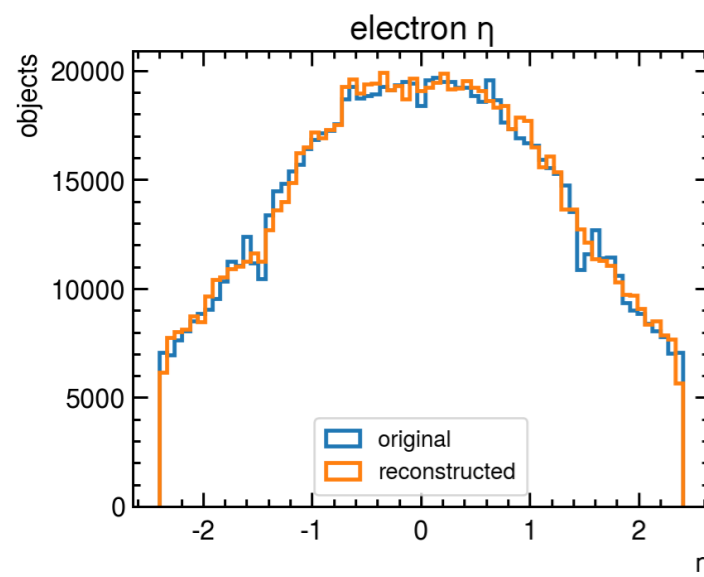
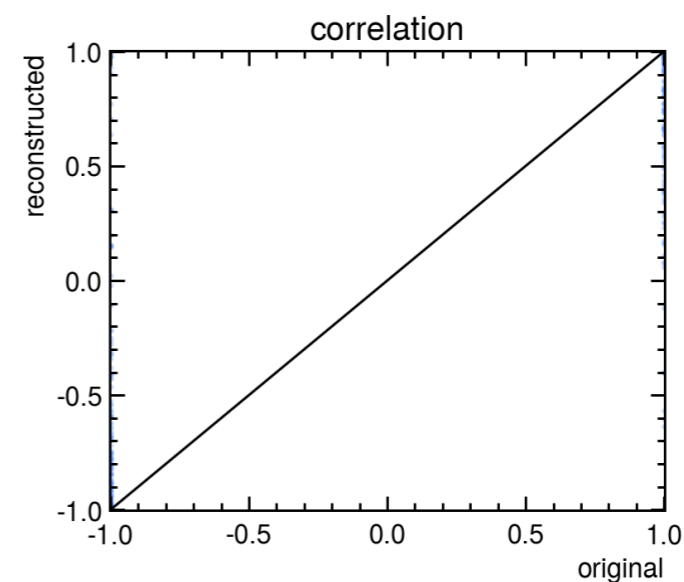
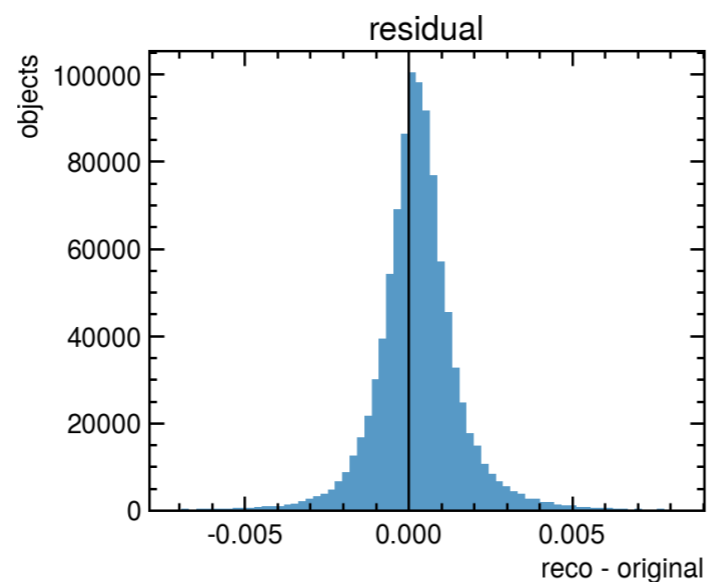
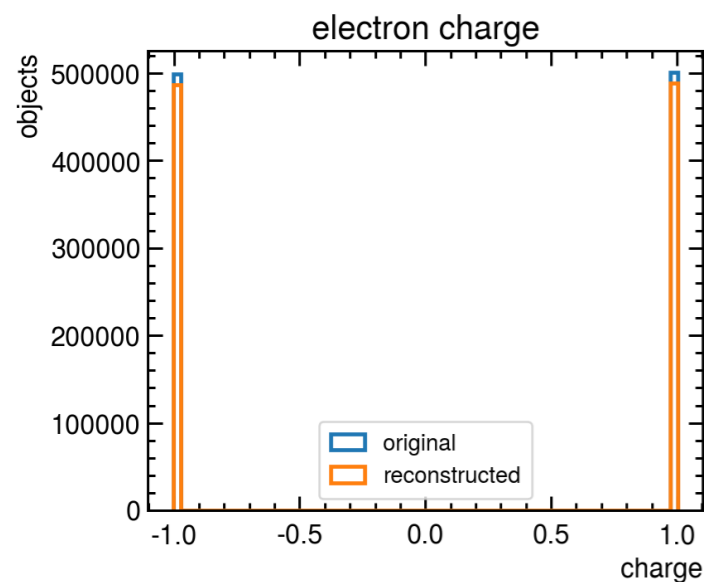
27.5.26

- * Last week, we agreed to create a PR for event-level tokenisation to 'heptokens' repo.
- * Link to PR: <https://github.com/Treasure-AmSC/heptokens/pull/5>
- * PR adds ATLAS event-object dataloading for per-object VQ-VAE tokenizer training (jets, electrons, muons, photons, taus), with configurable H5 feature paths and tokenizer diagnostics.
- * Ran over the recent files from Viviana, next slides for the results

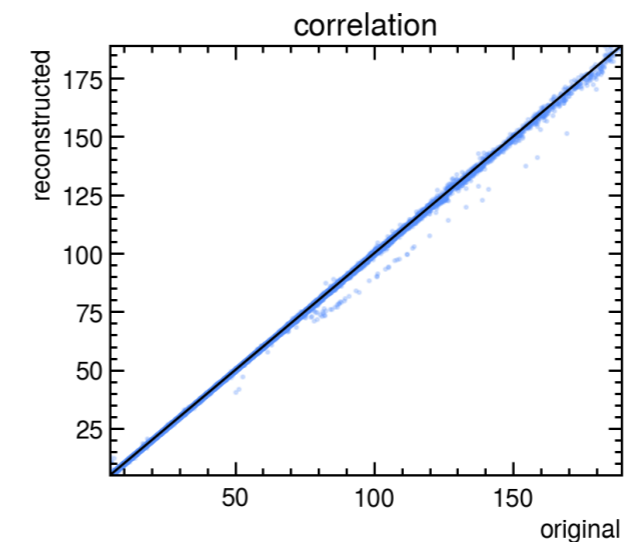
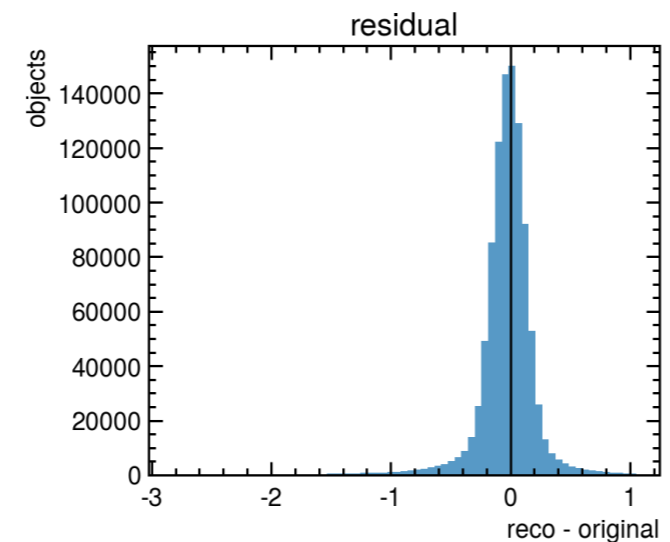
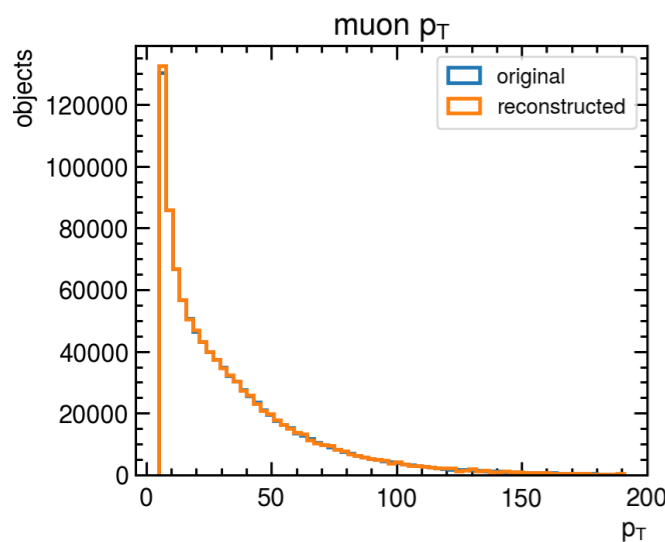
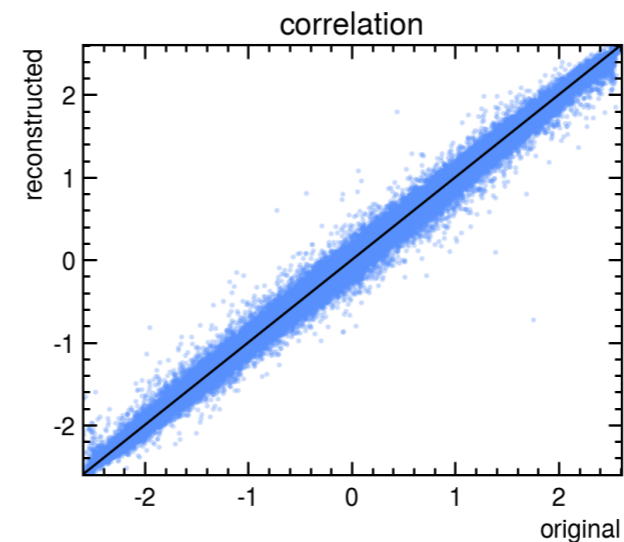
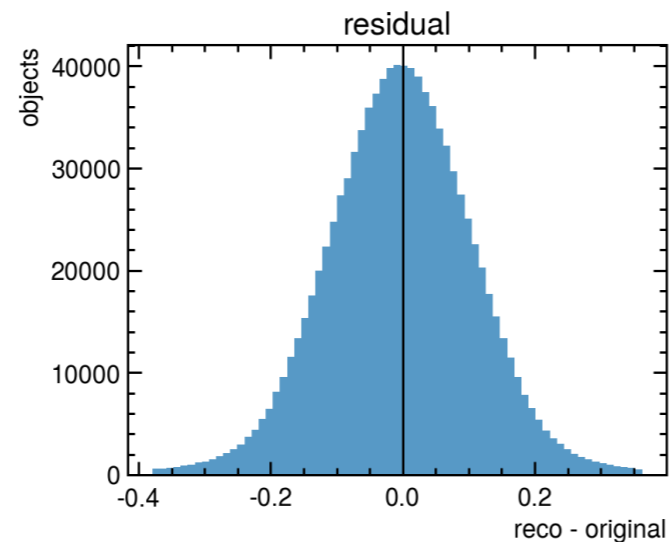
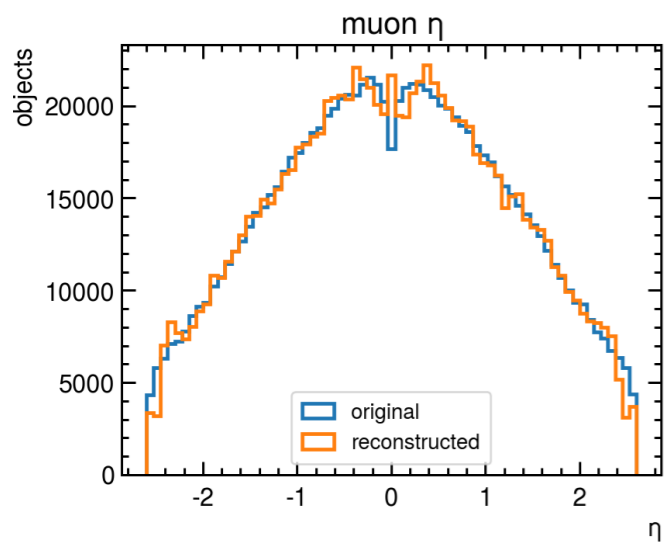
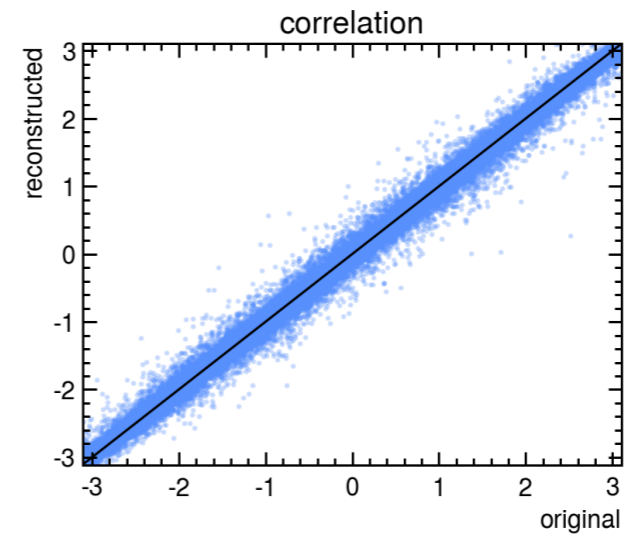
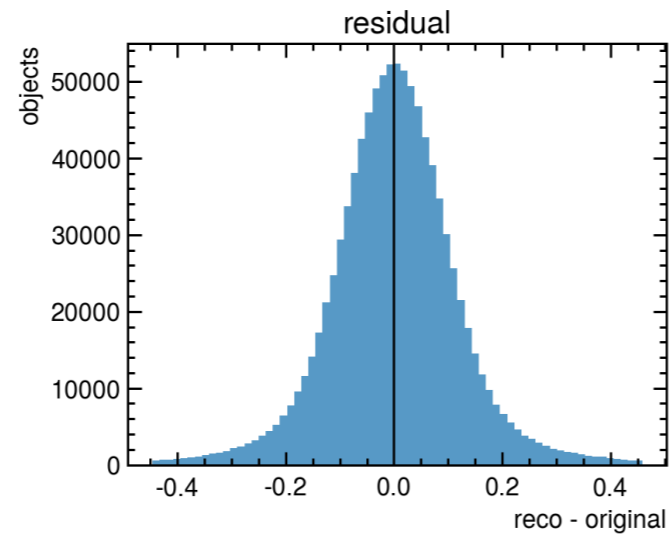
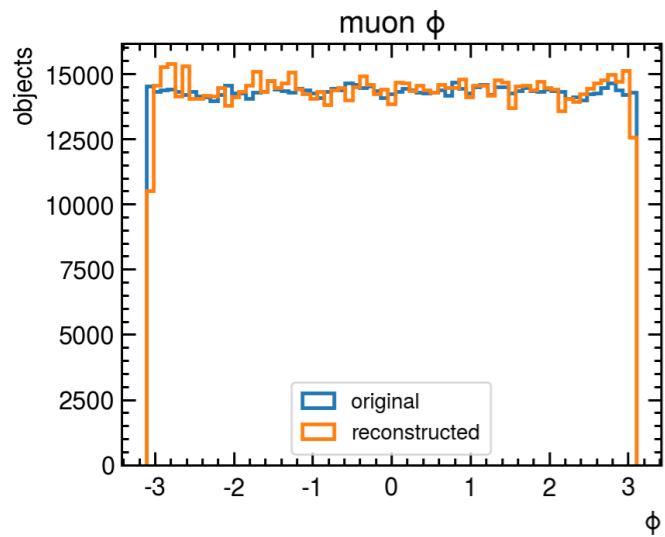
```
* jets      events=3,558,000 objects=14,344,784
  electrons events=3,558,000 objects=3,150,807
  muons     events=3,558,000 objects=3,945,218
  photons   events=3,558,000 objects=688,912
  taus      events=3,558,000 objects=4,260,772
```



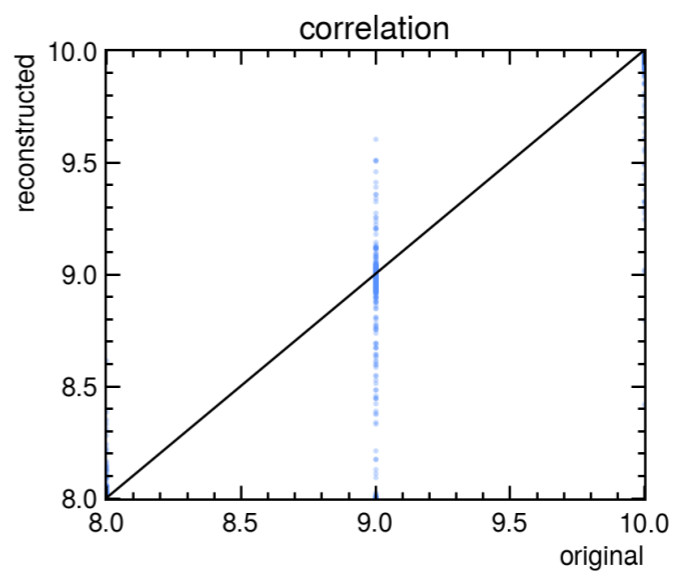
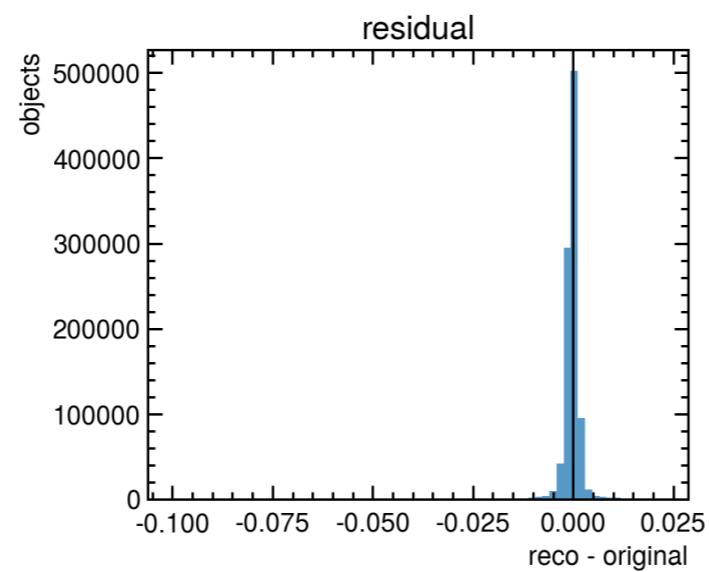
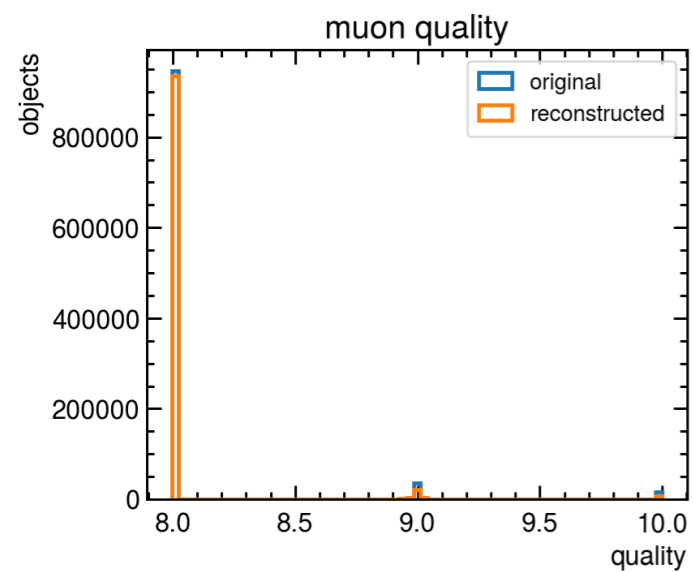
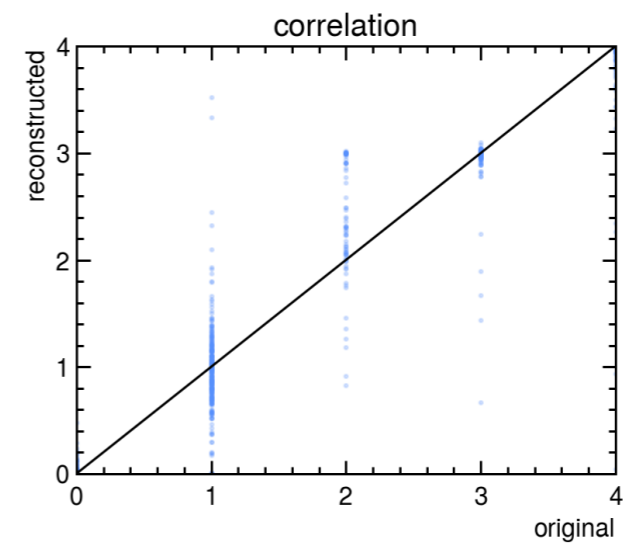
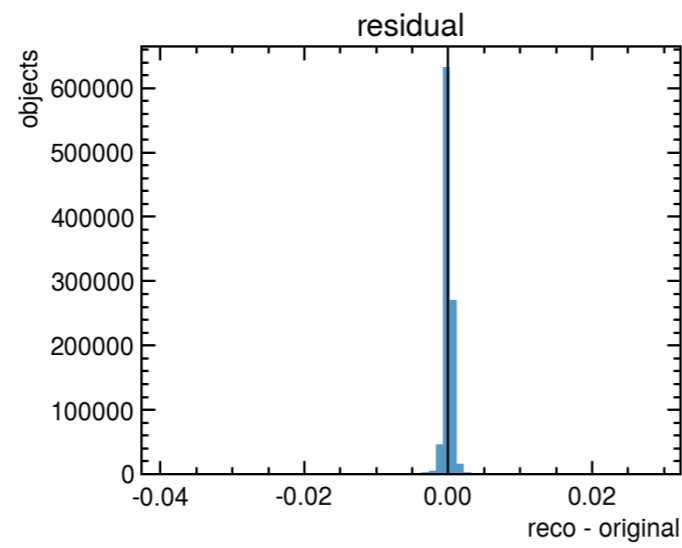
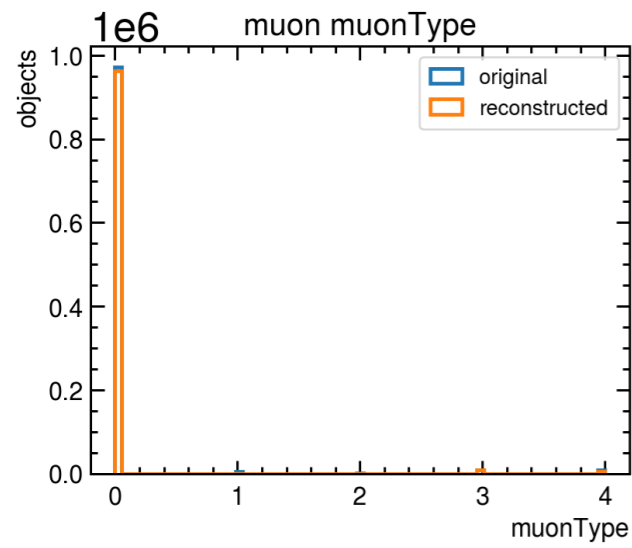
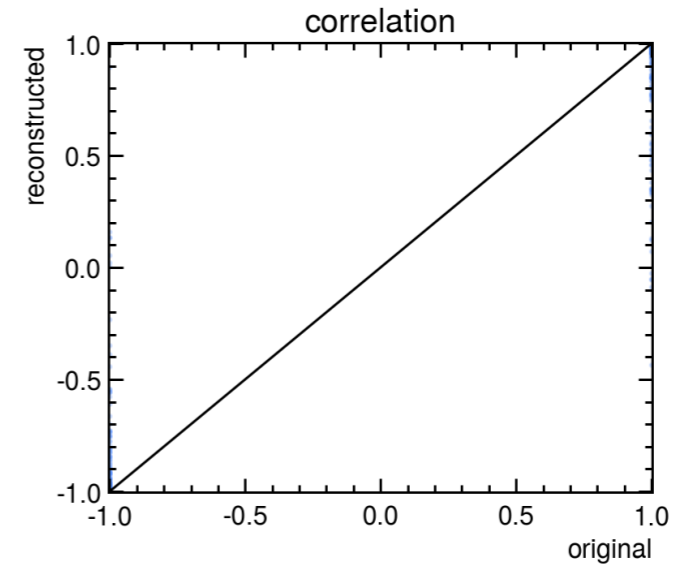
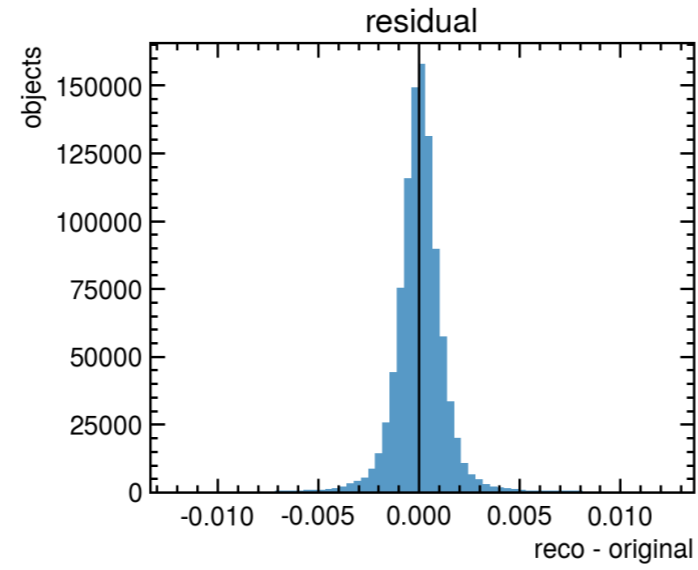
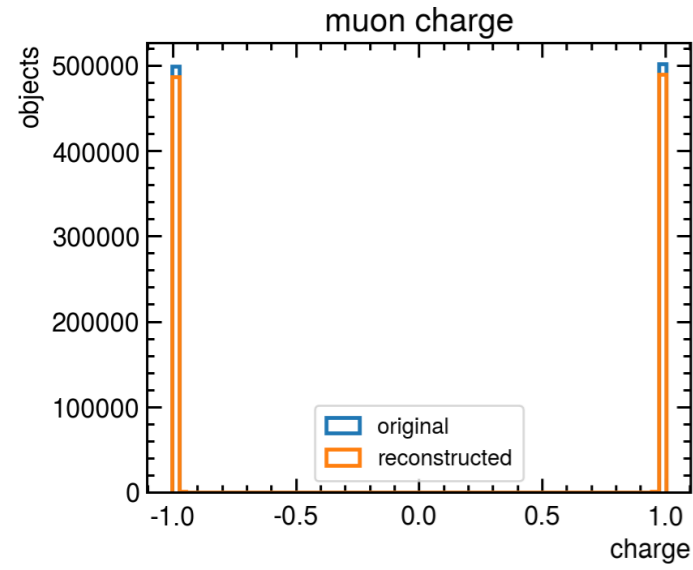
Diagnostics - electrons



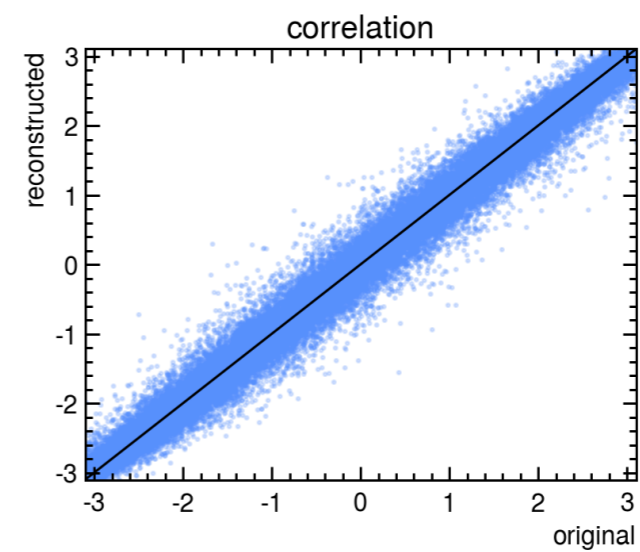
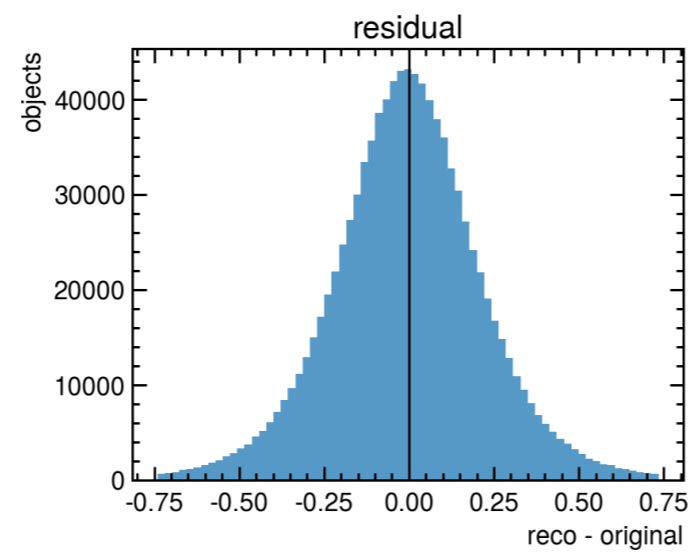
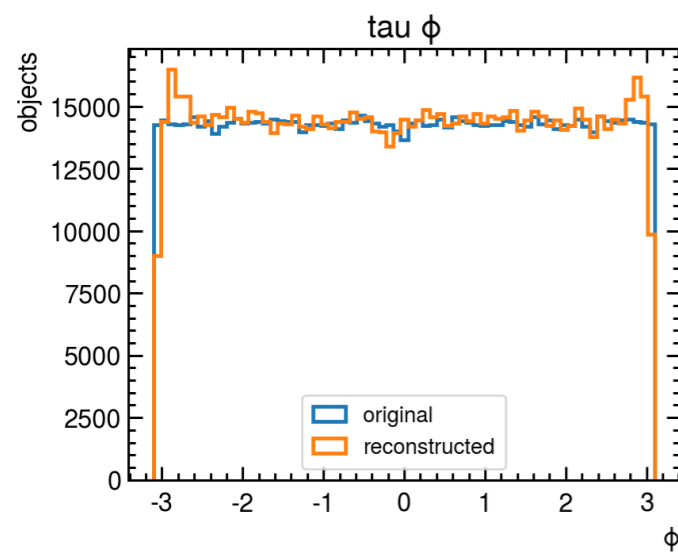
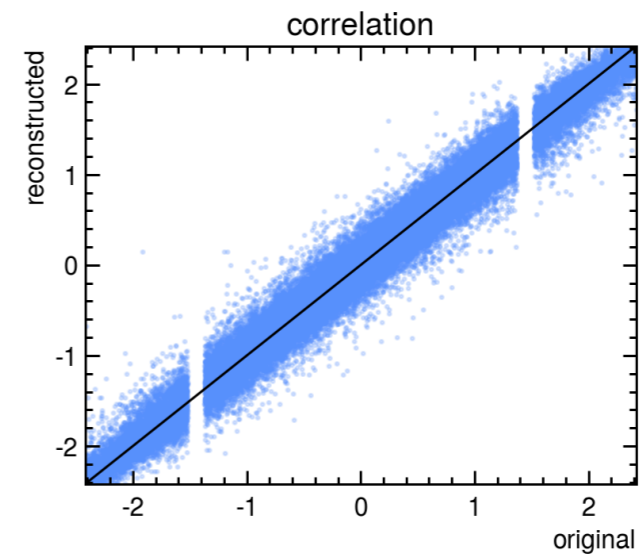
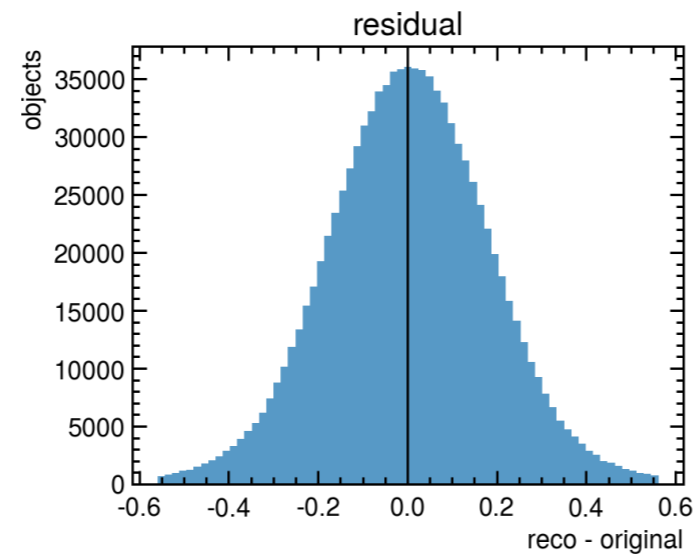
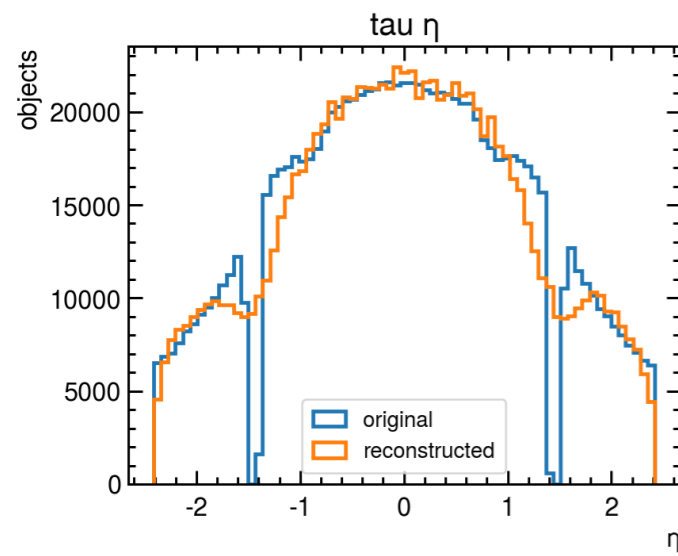
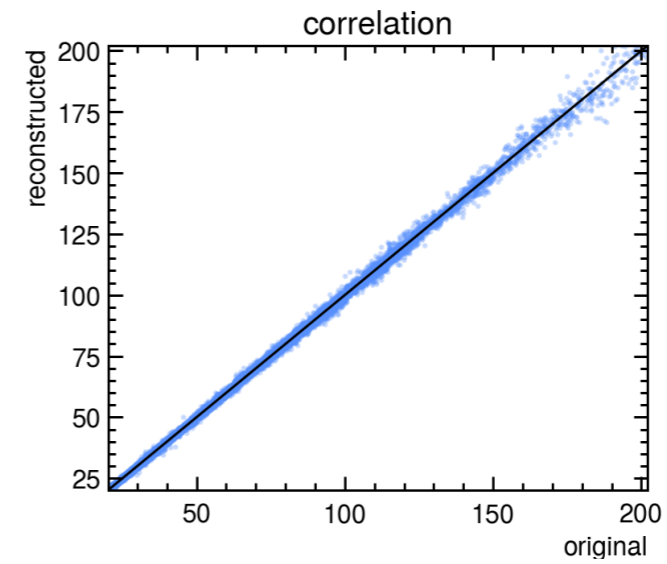
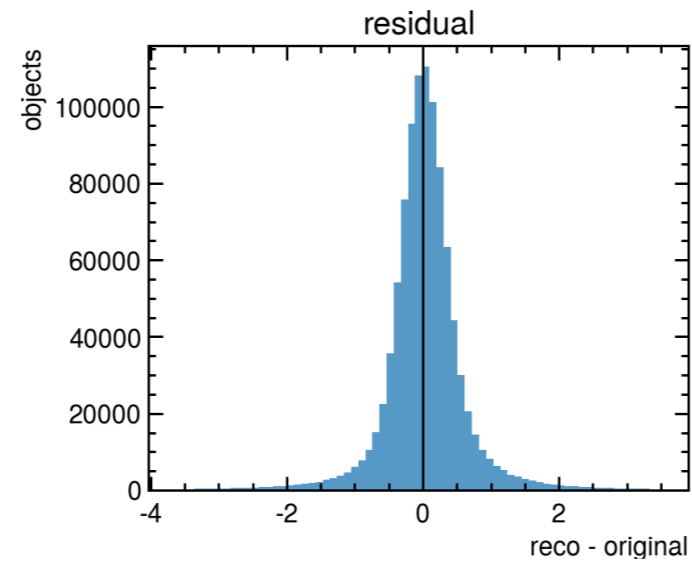
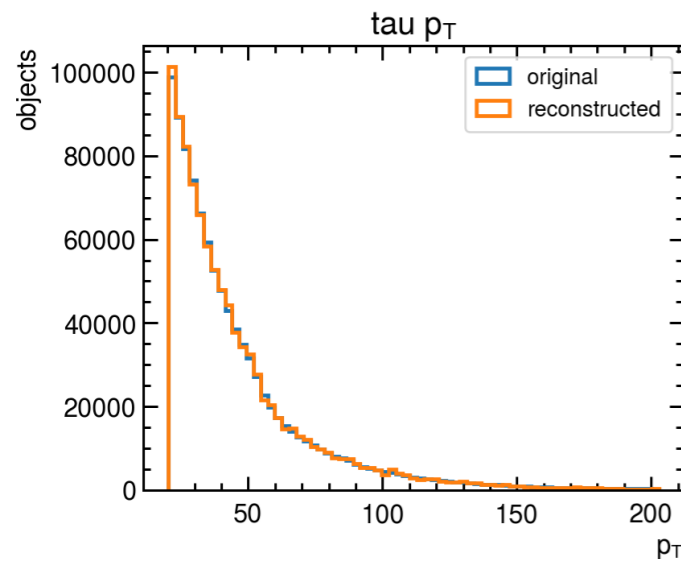
Diagnostics - muons



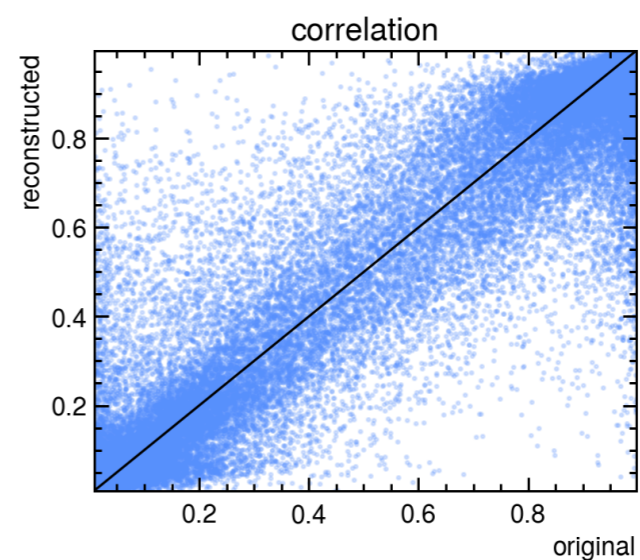
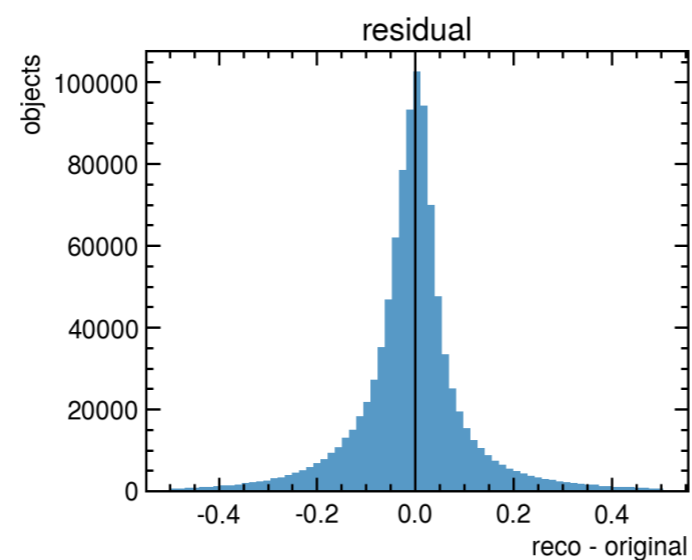
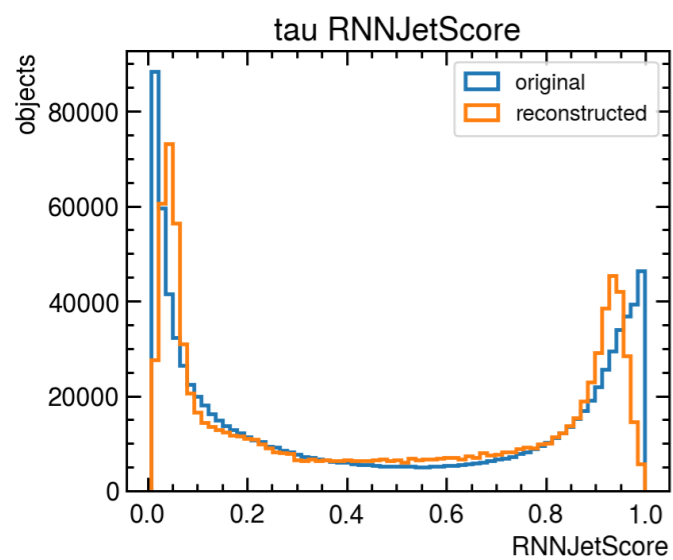
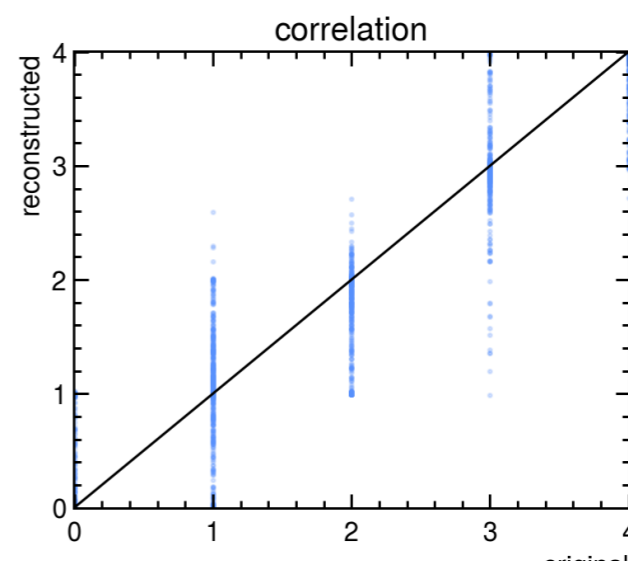
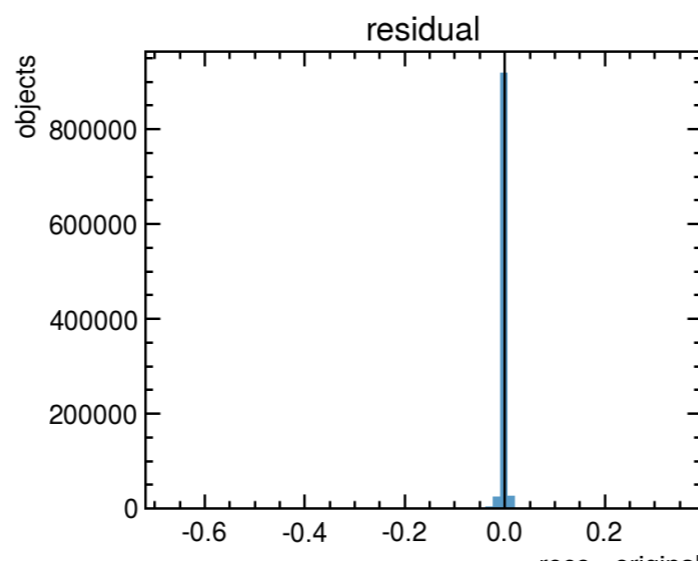
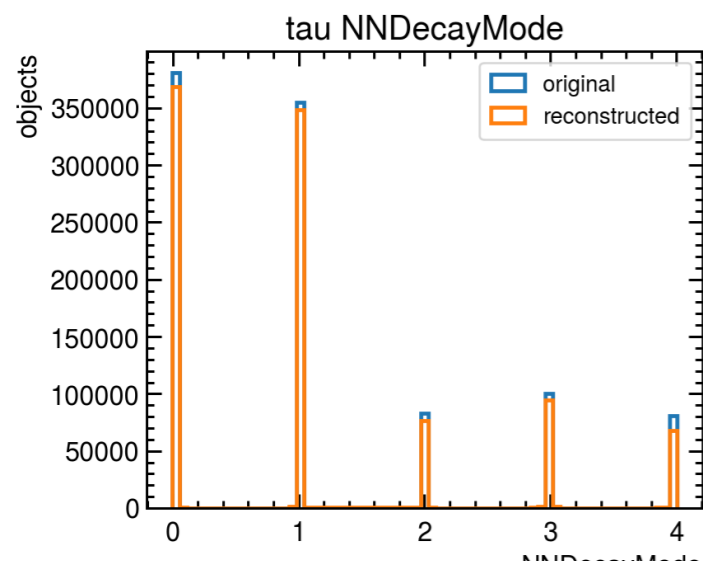
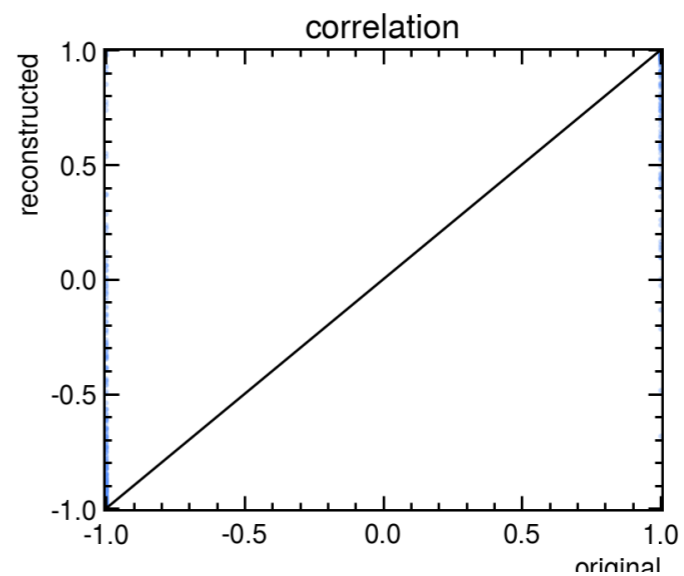
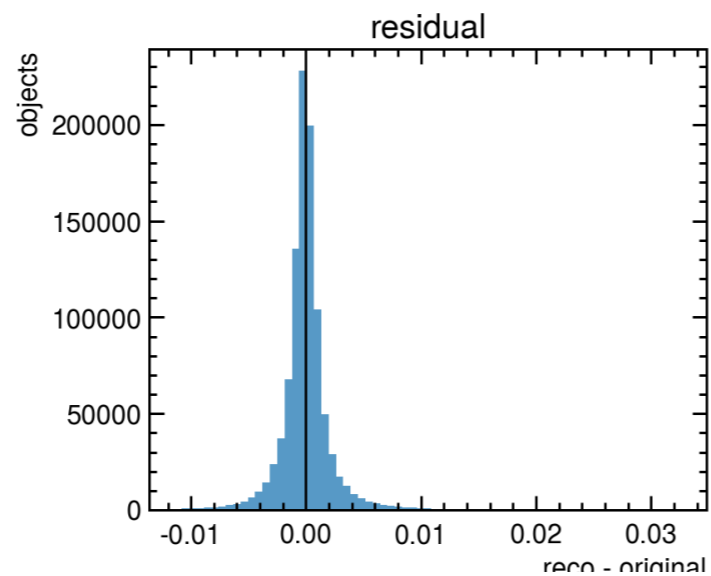
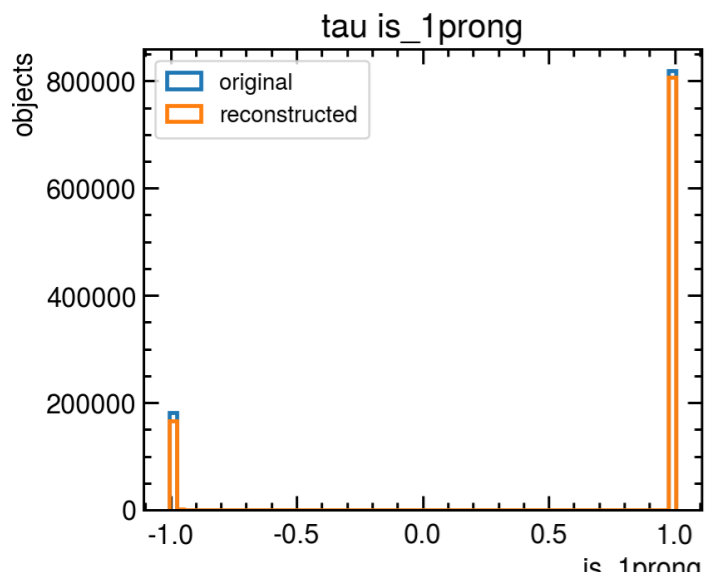
Diagnostics - muons



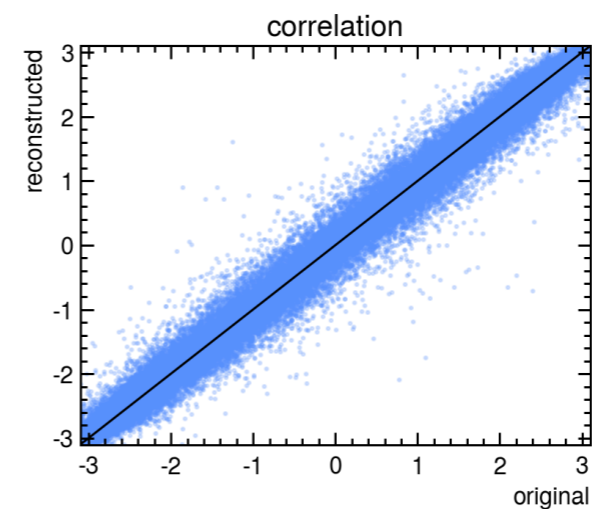
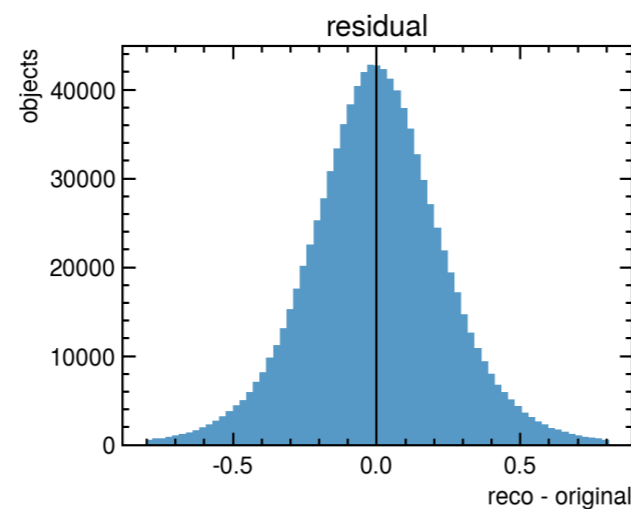
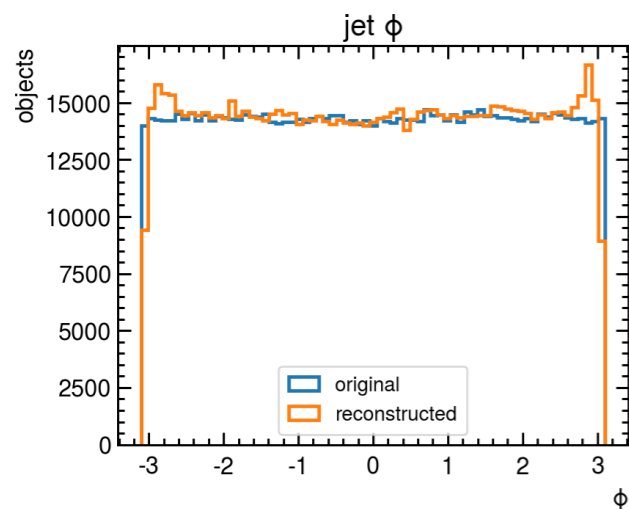
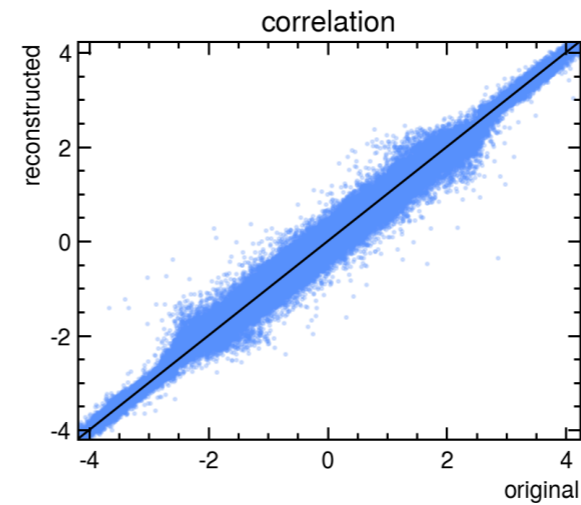
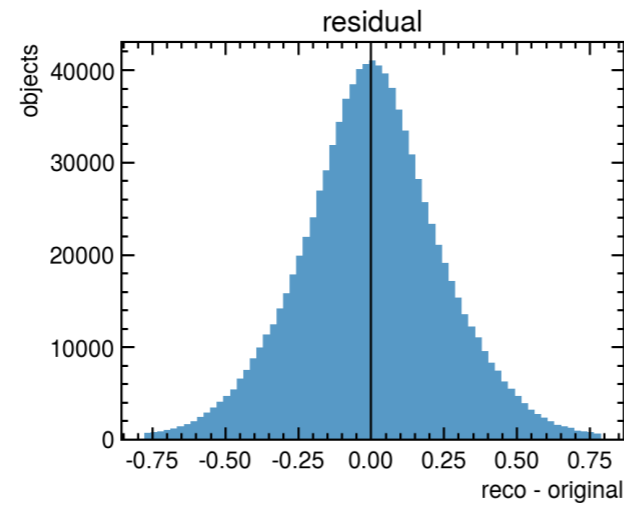
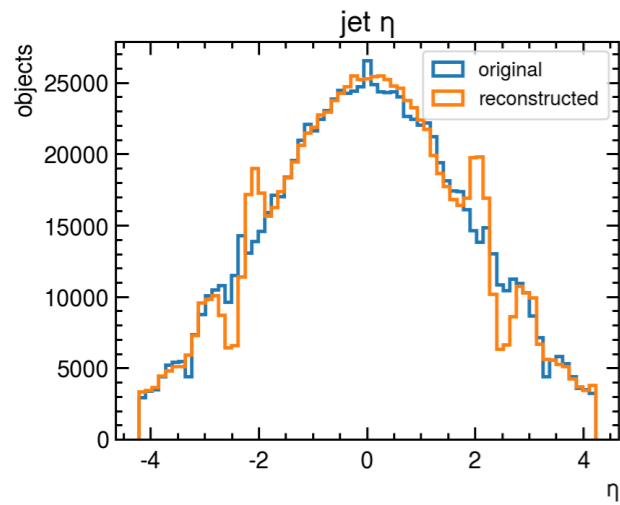
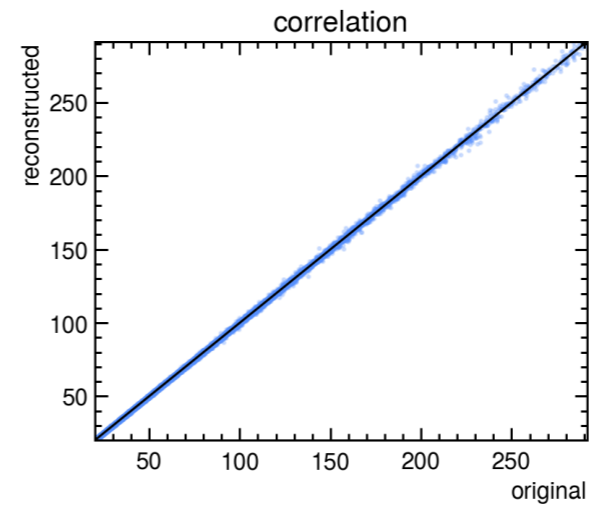
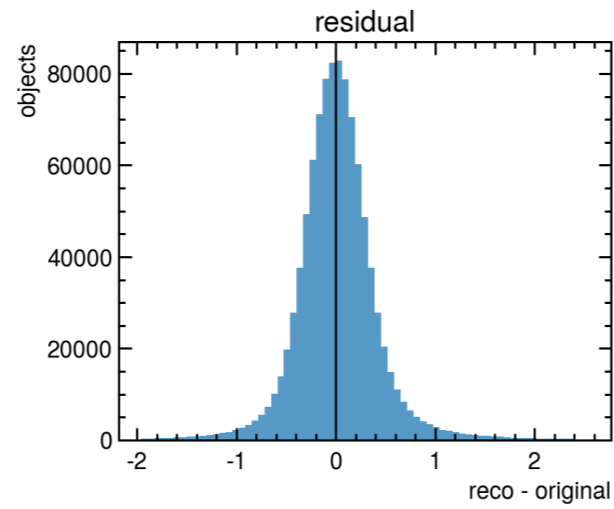
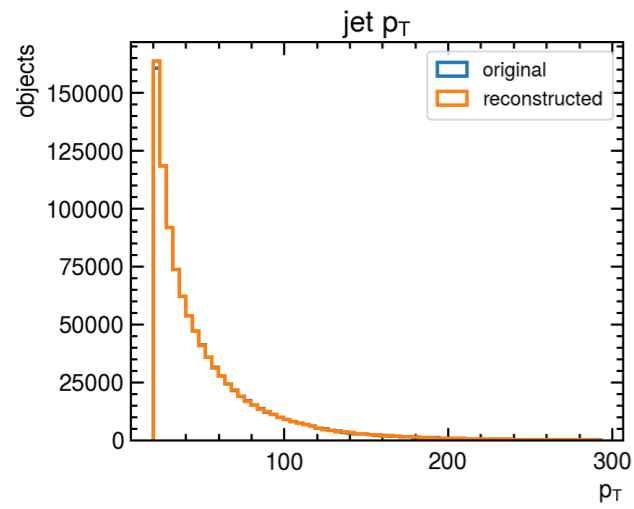
Diagnostics - taus



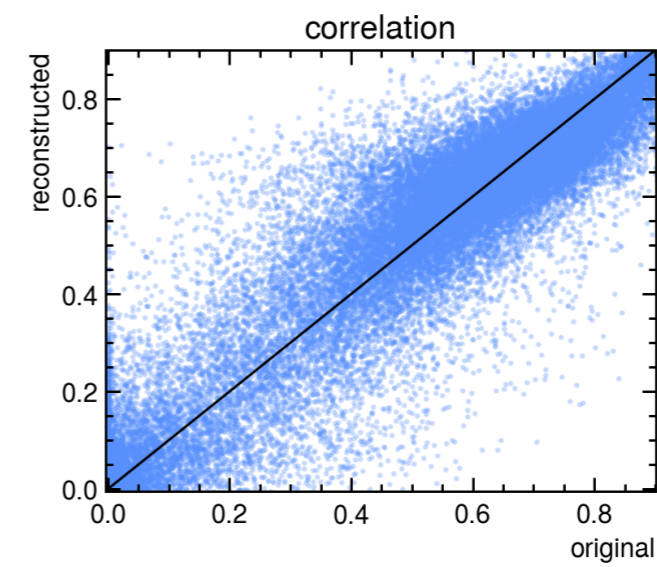
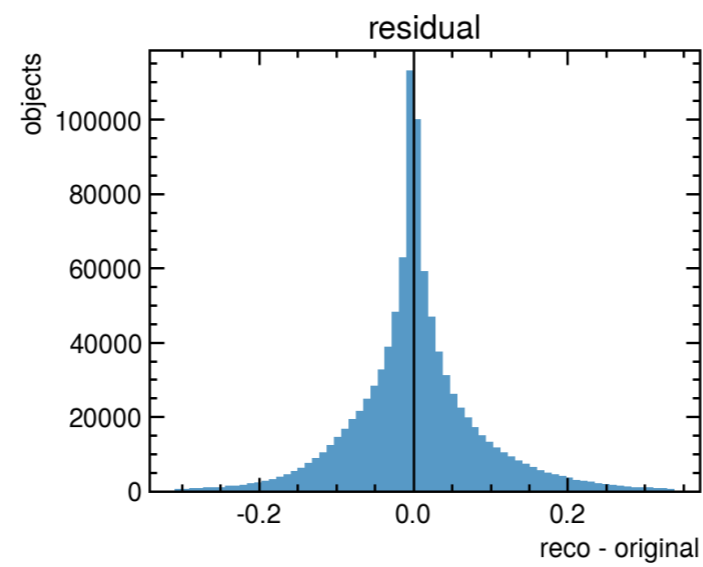
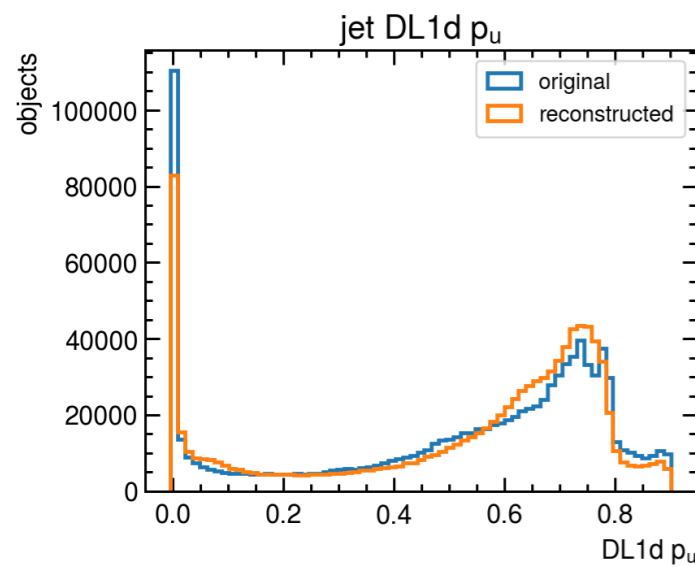
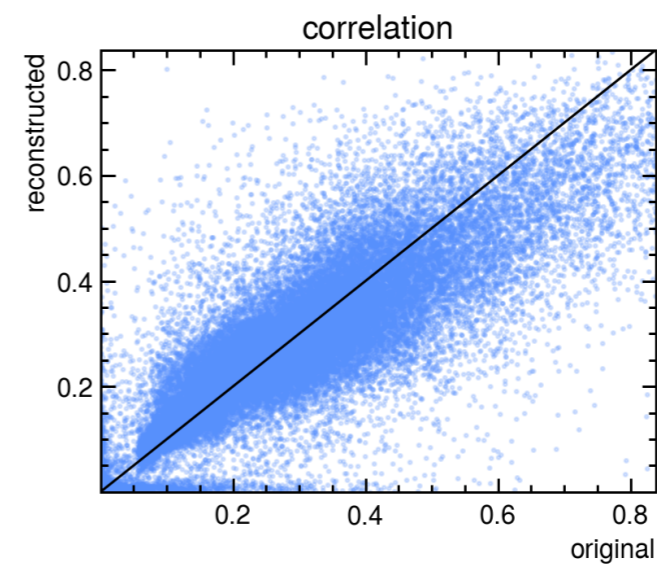
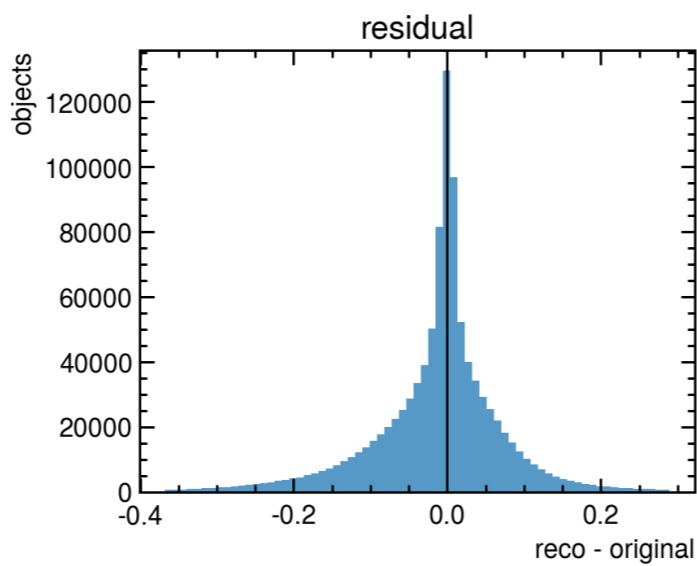
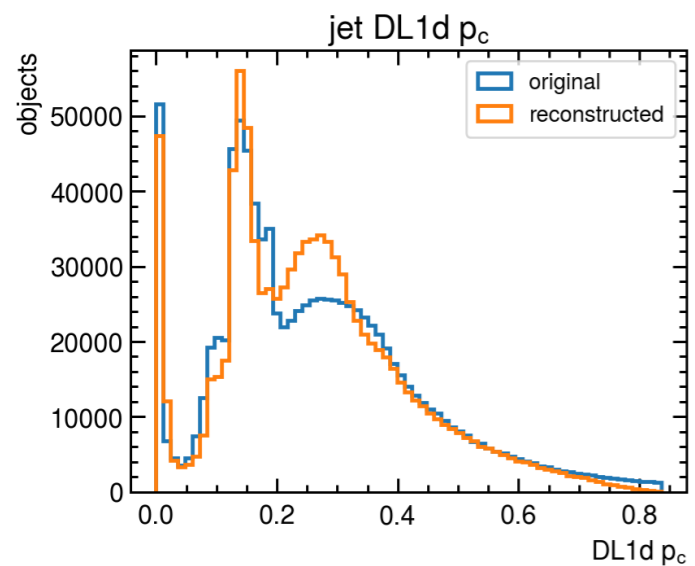
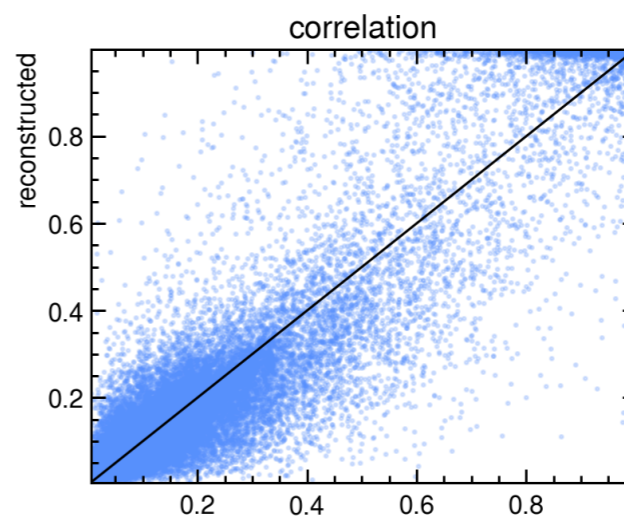
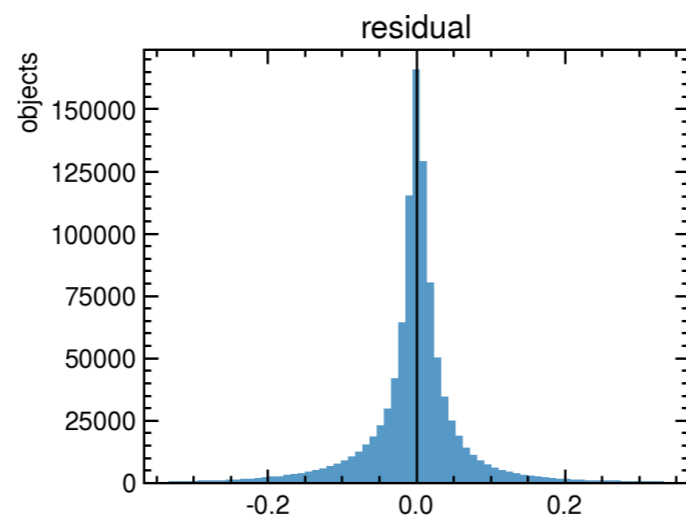
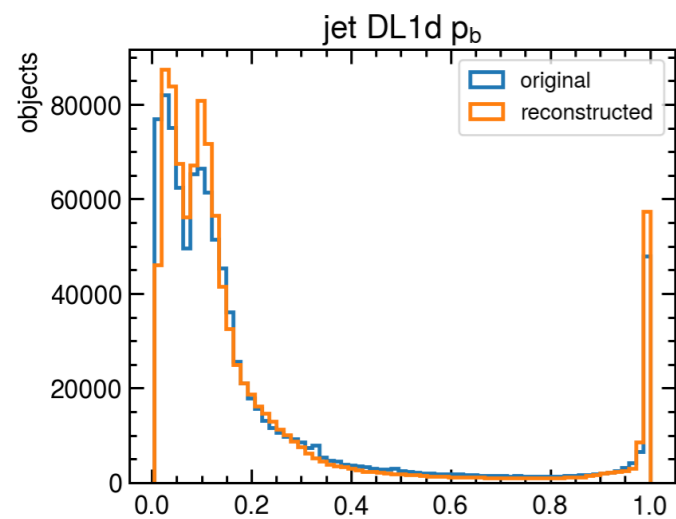
Diagnostics - taus



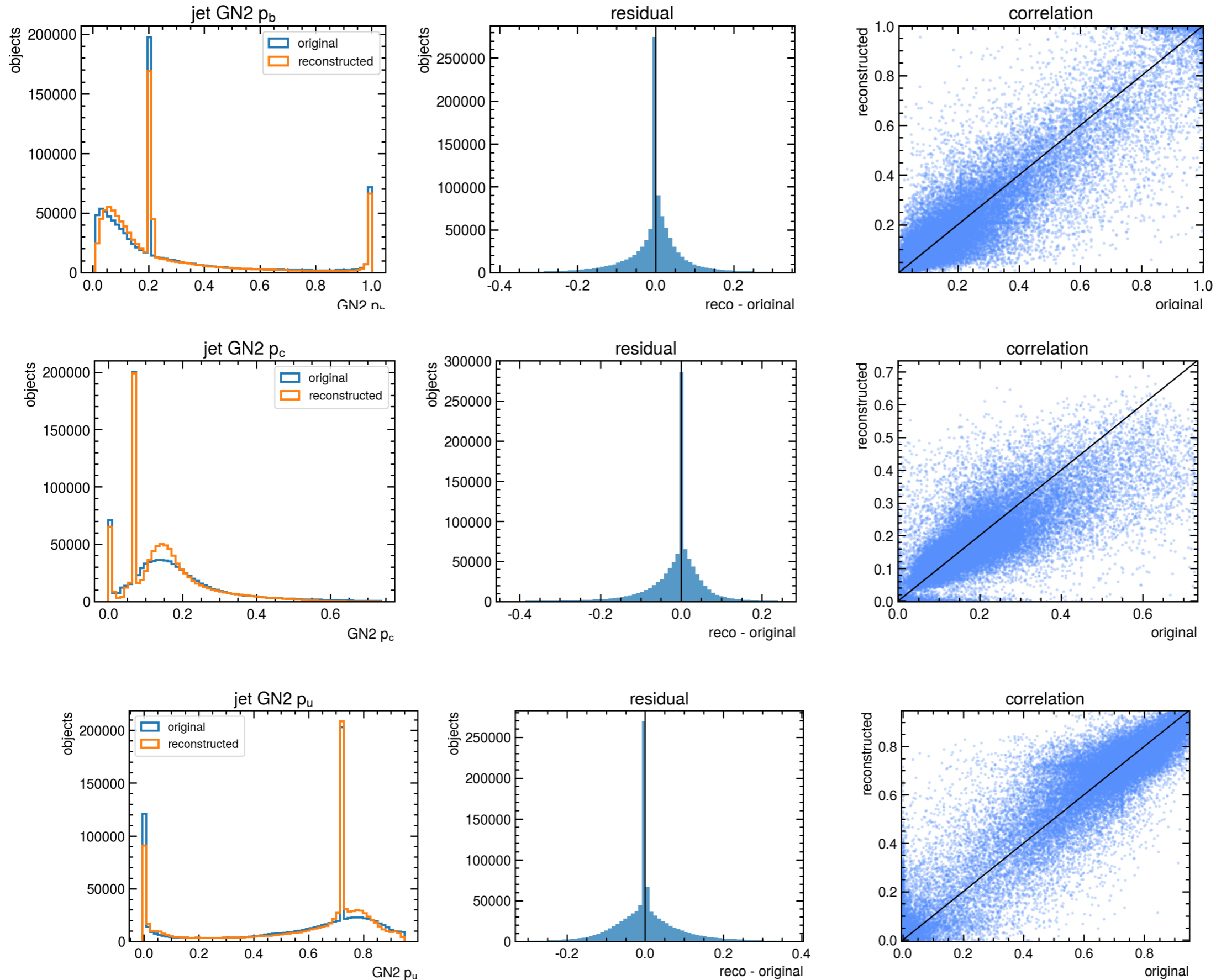
Diagnostics - jets



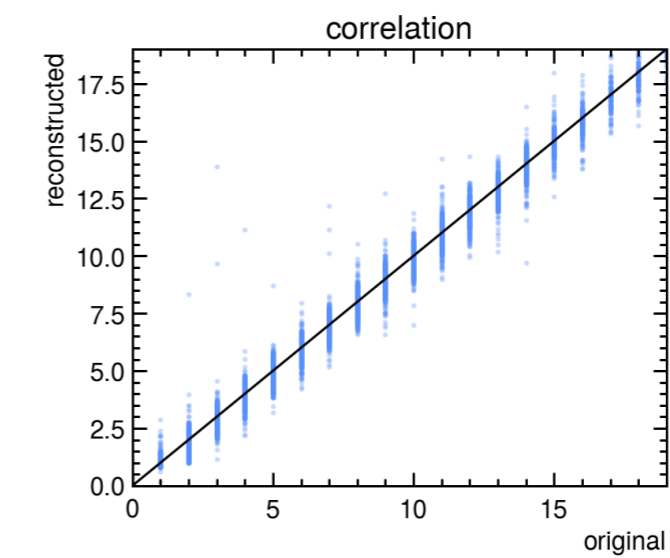
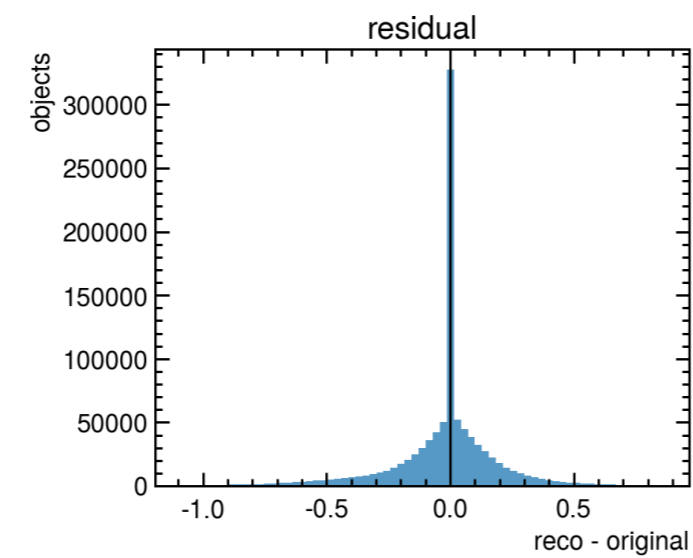
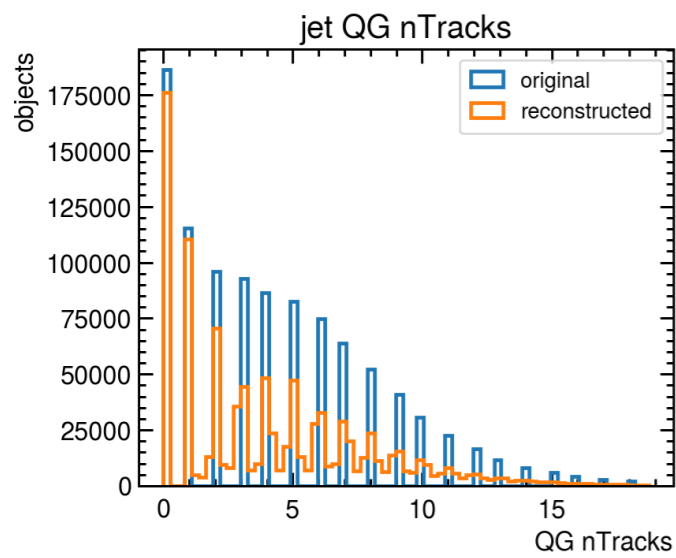
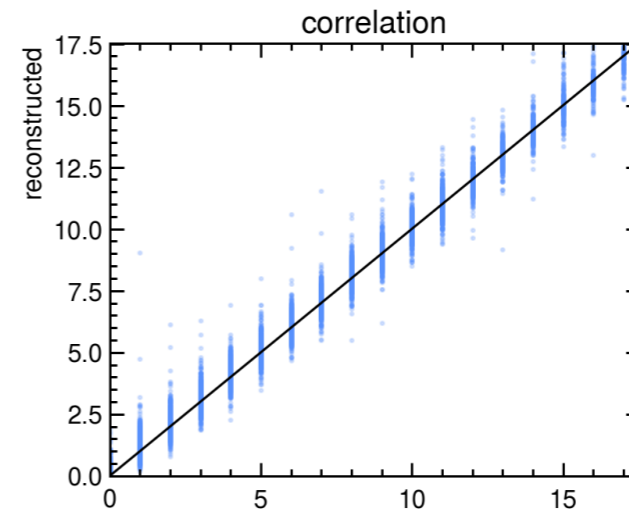
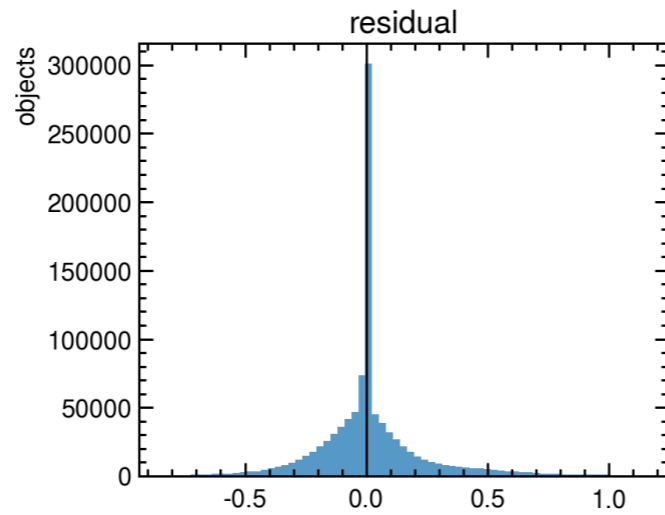
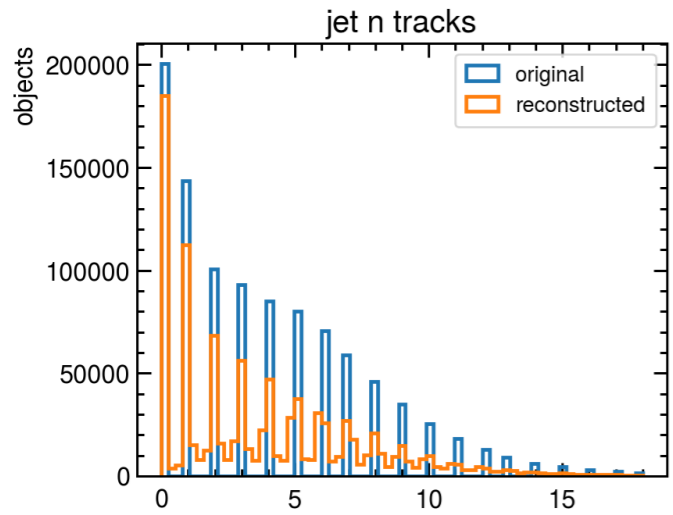
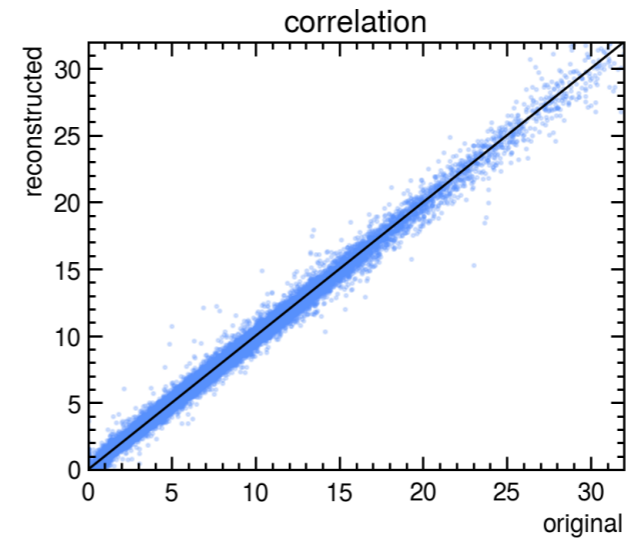
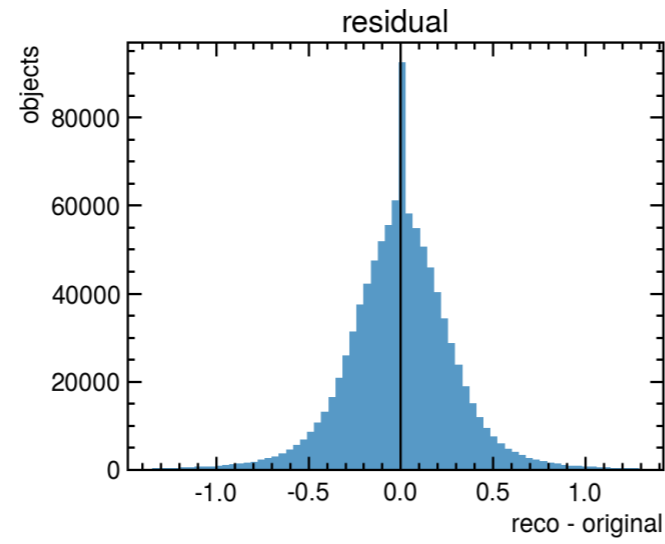
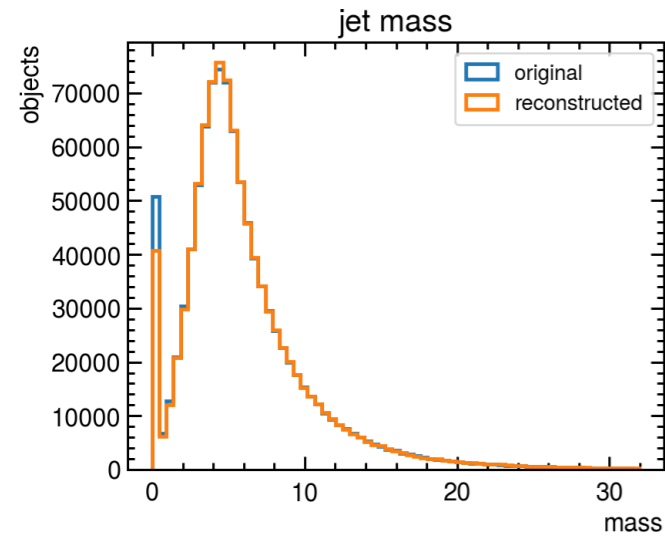
Diagnostics - jets



Diagnostics - jets

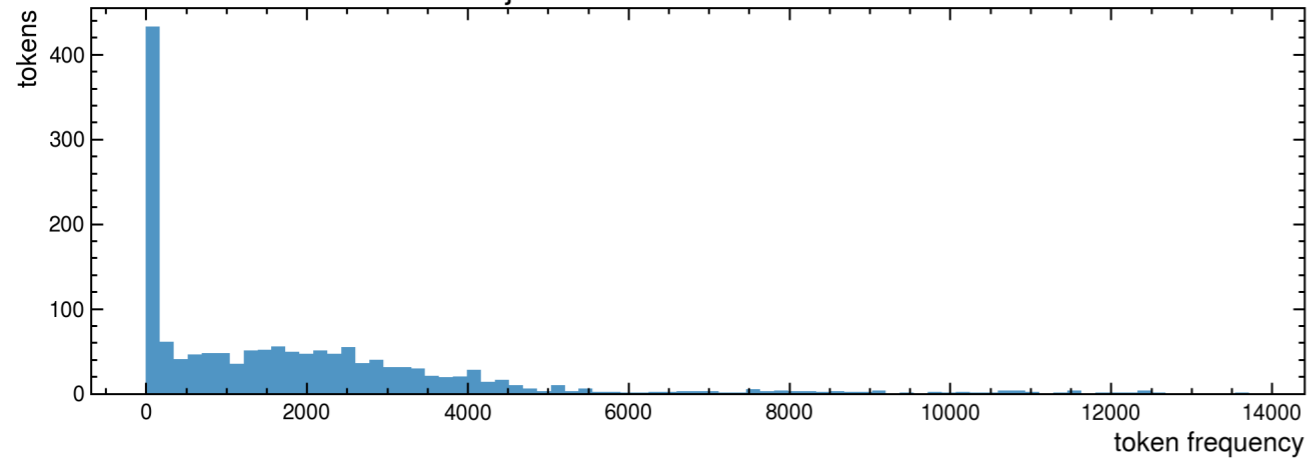


Diagnostics - jets



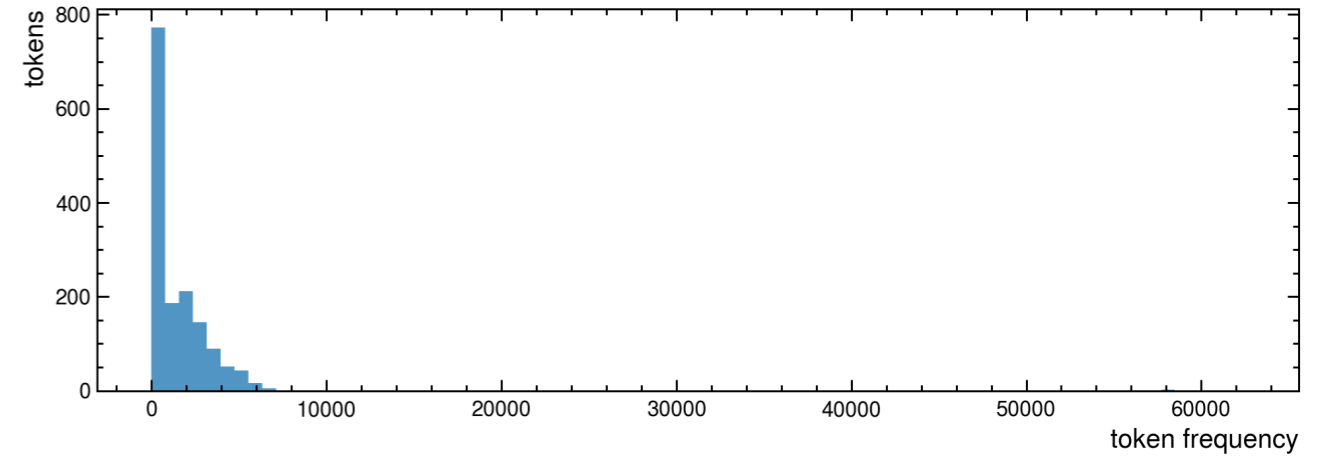
Diagnostics

jet codebook: 75.3% used



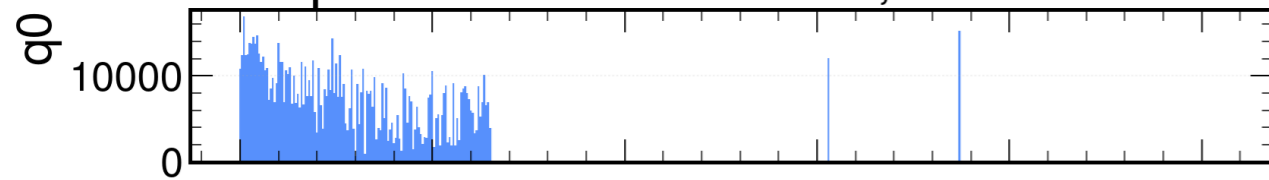
Codebook usage frequency

electron codebook: 67.8% used

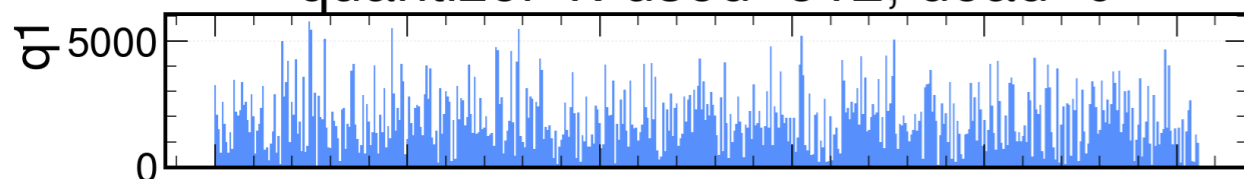


Codebook usage frequency

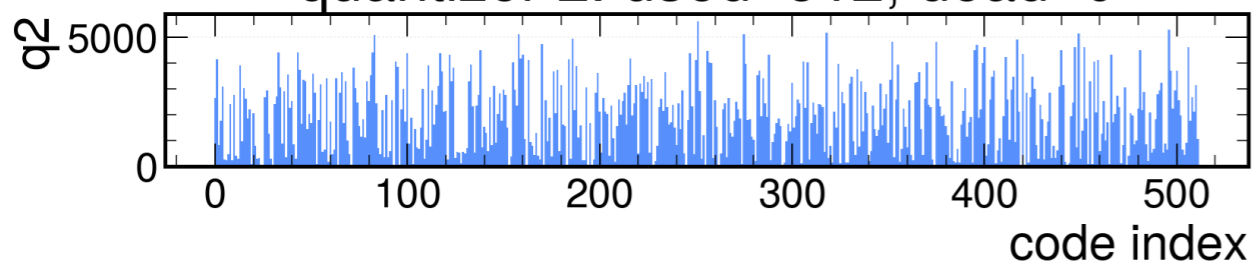
quantizer 0: used=133, dead=379



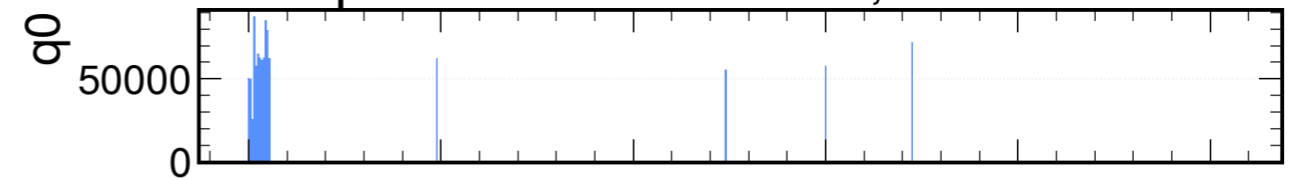
quantizer 1: used=512, dead=0



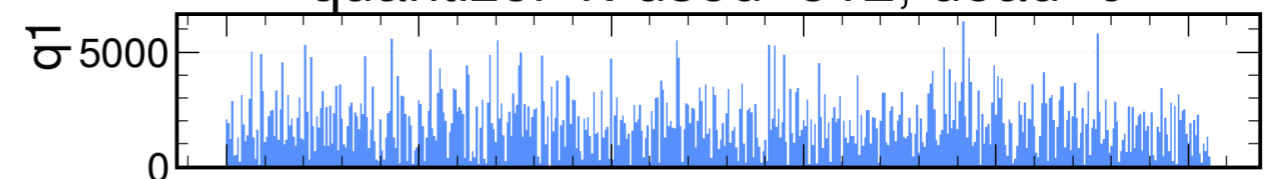
quantizer 2: used=512, dead=0



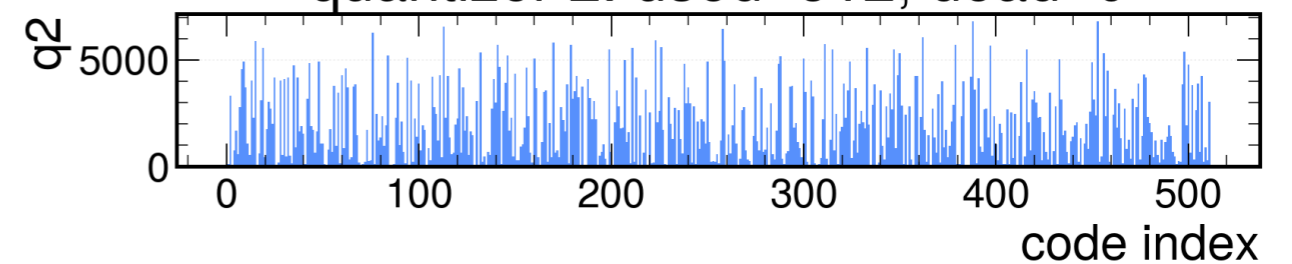
quantizer 0: used=17, dead=495



quantizer 1: used=512, dead=0

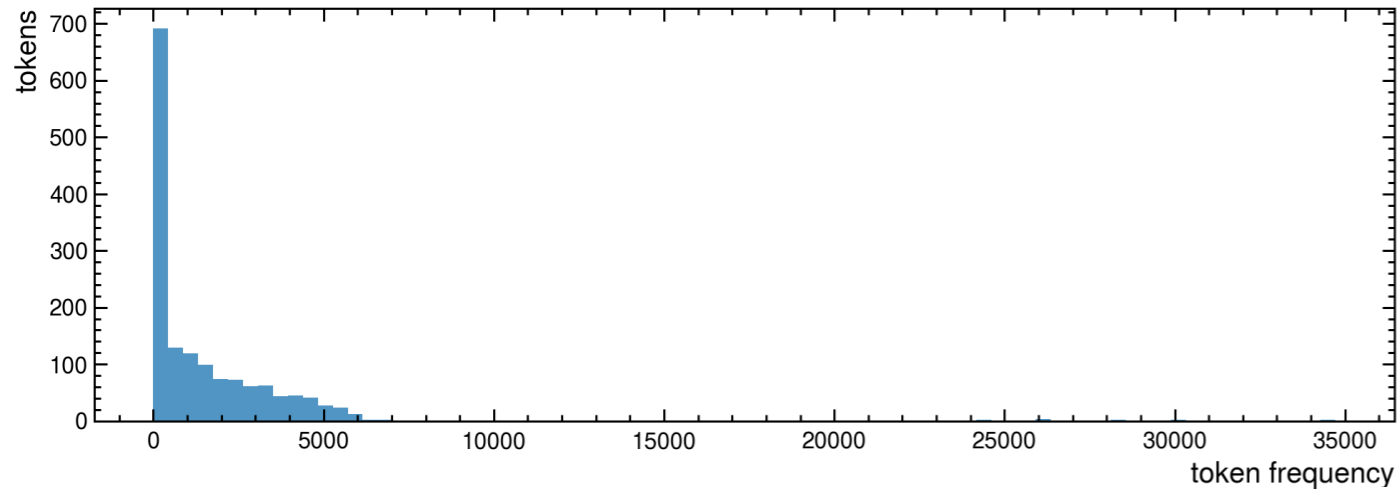


quantizer 2: used=512, dead=0

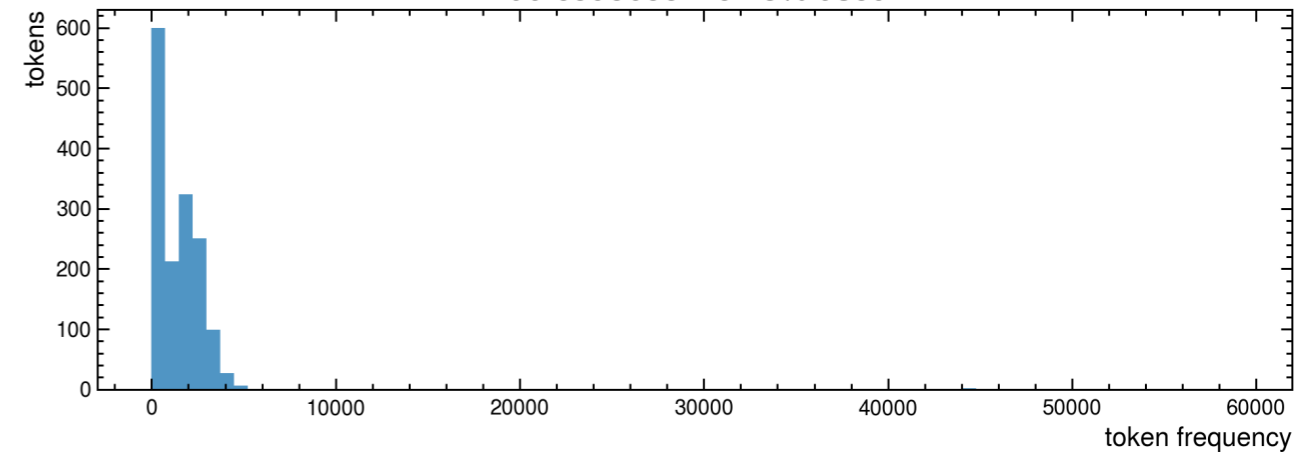


Diagnostics

muon codebook: 68.6% used

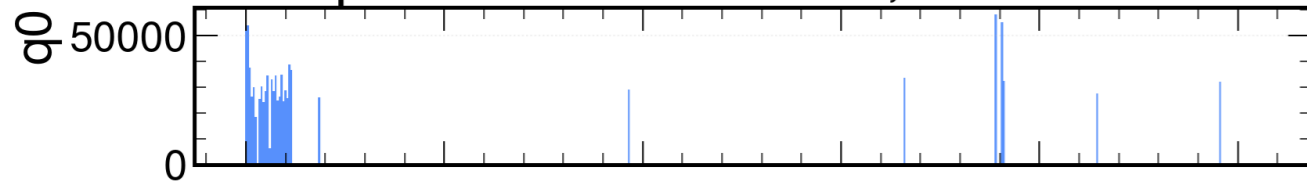


tau codebook: 67.8% used

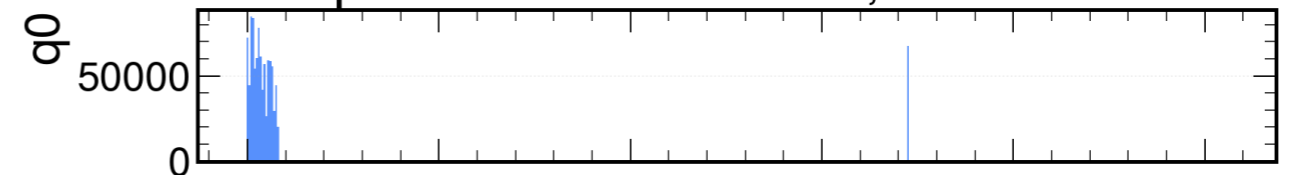


Codebook usage frequency

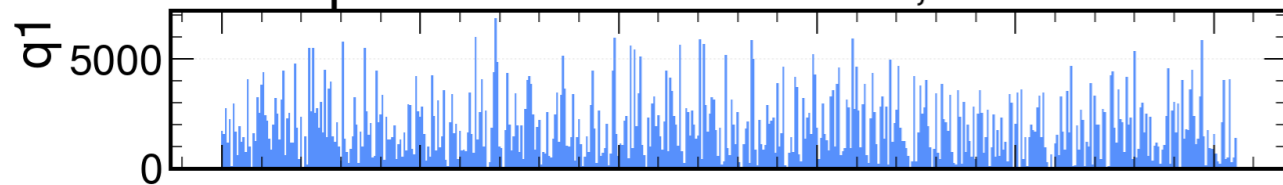
quantizer 0: used=32, dead=480



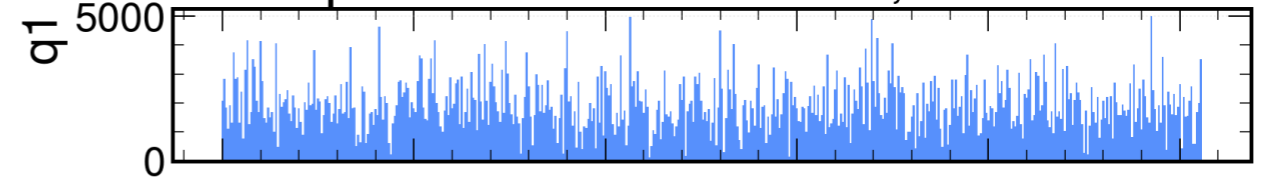
quantizer 0: used=18, dead=494



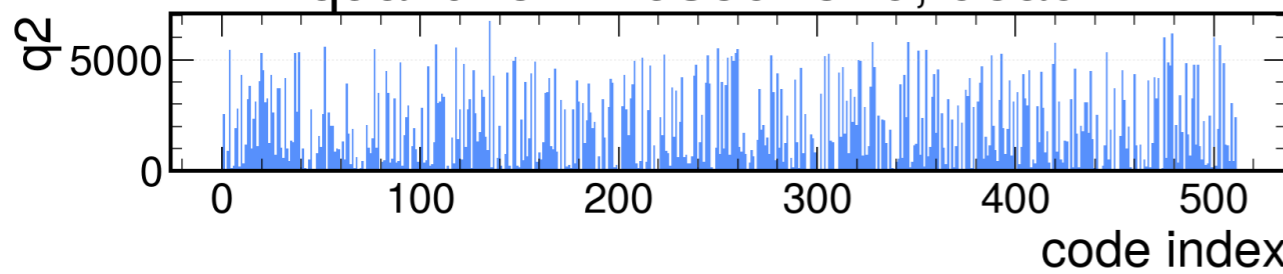
quantizer 1: used=511, dead=1



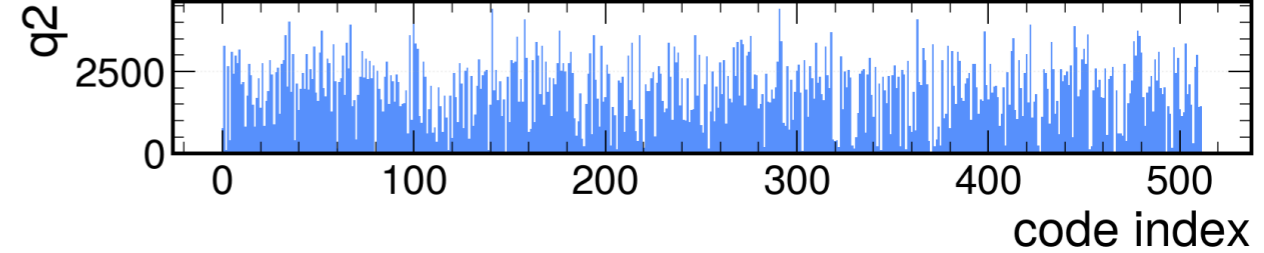
quantizer 1: used=512, dead=0



quantizer 2: used=510, dead=2



quantizer 2: used=512, dead=0



Discussion

- * First full-feature object tokenizers are now trained and diagnostics are running for all object types.
- * Basic kinematics are reconstructed well, but some higher-level/discriminant features are not yet well preserved.
 - ▶ The VQ-VAE appears to smooth bounded score-like features such as DL1d/GN2 probabilities.
- * Next steps: Add data, separate continuous, discrete, and bounded score-like variables more carefully; test feature-wise preprocessing/normalisation; compare **single-codebook vs residual-codebook settings**; and evaluate whether these reconstruction differences matter for downstream model.