

Treasure event level-tokenisation

Merve Nazlim Agaras

3.6.26

Stats of the events

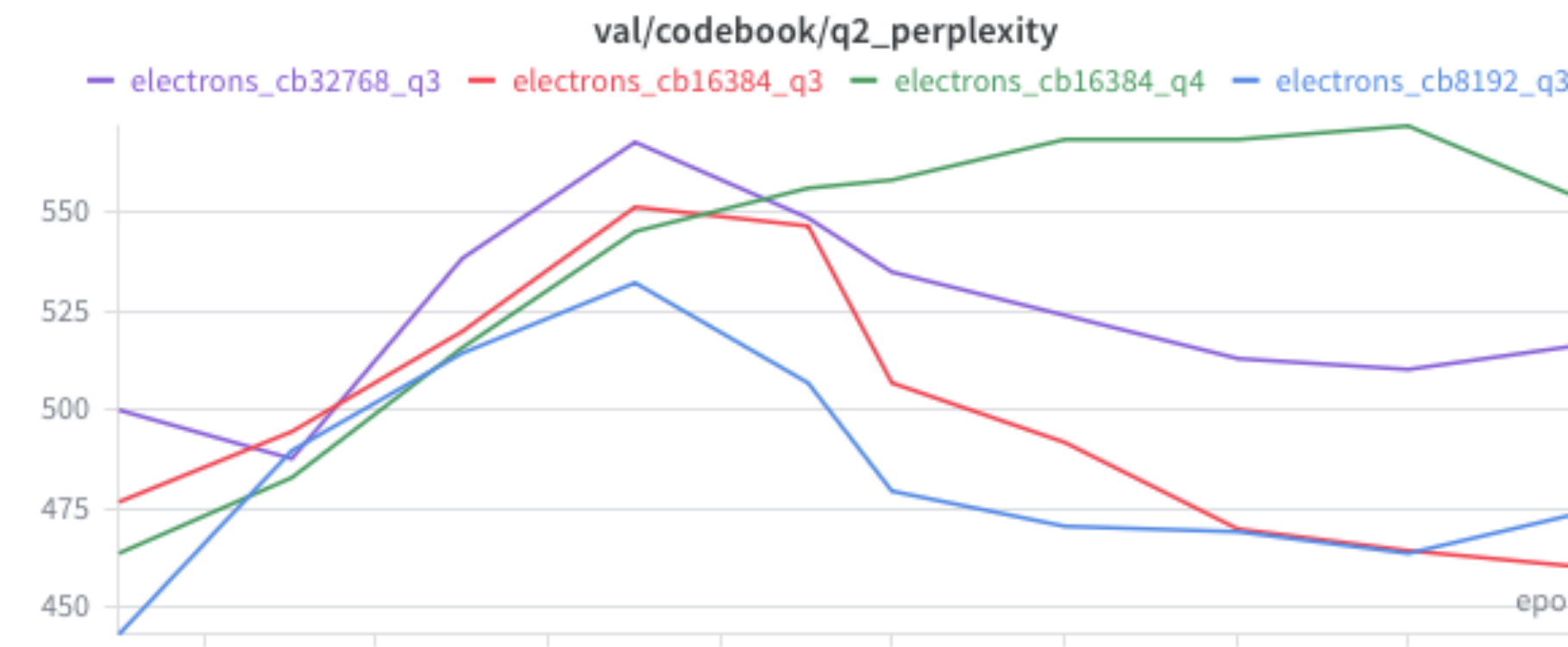
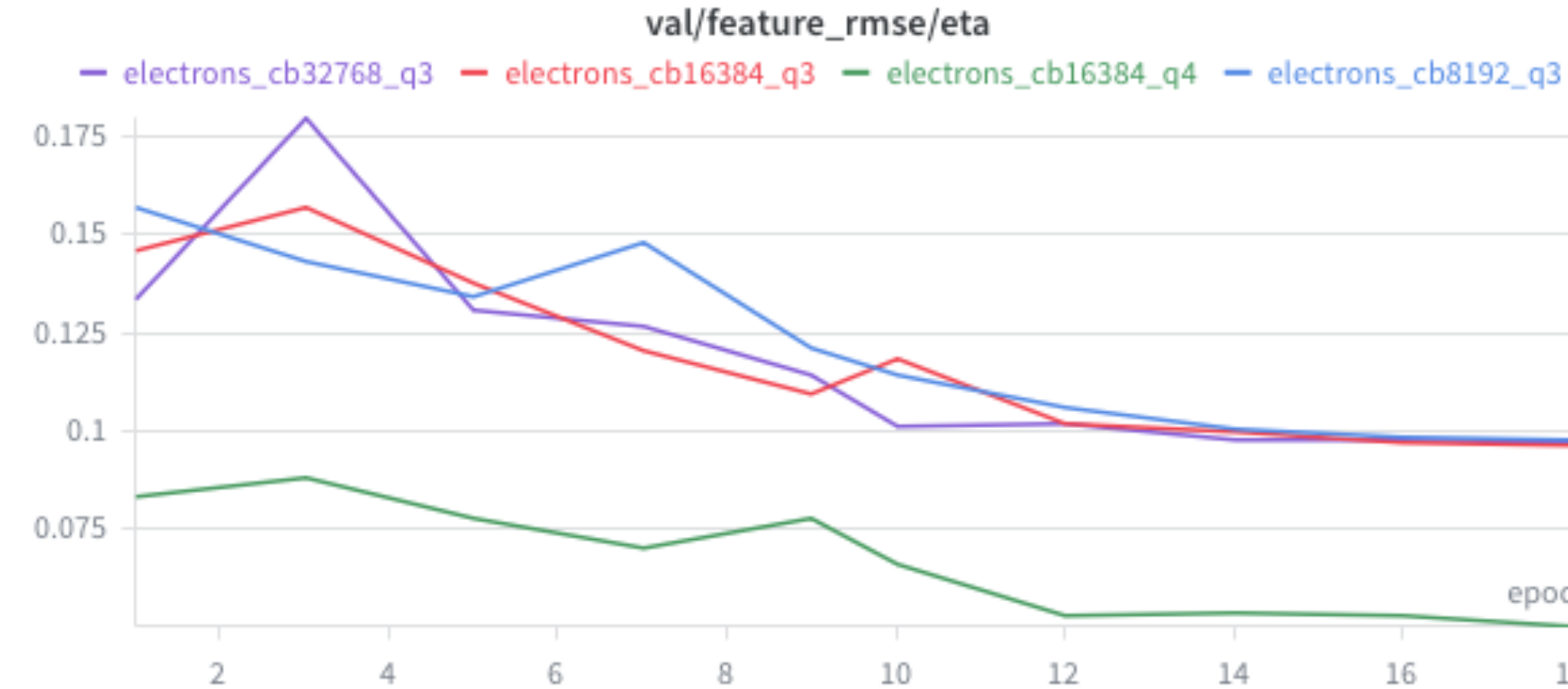
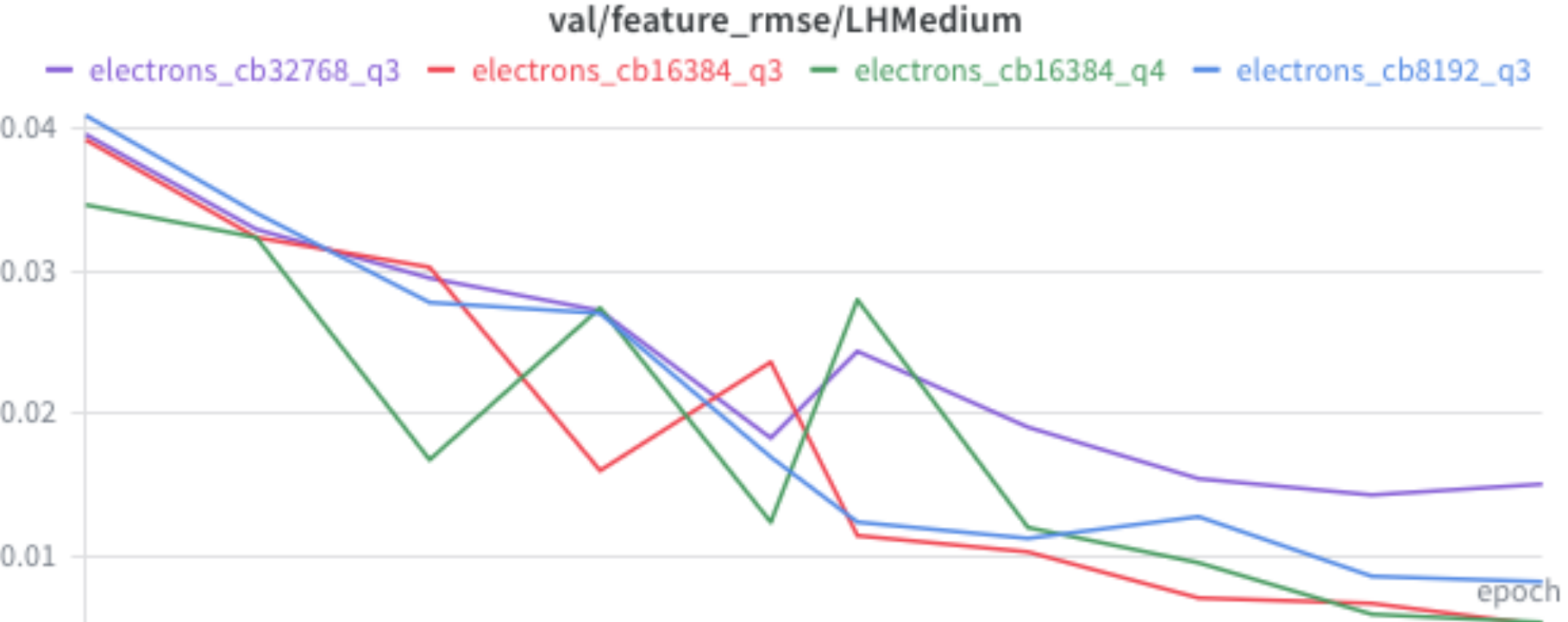
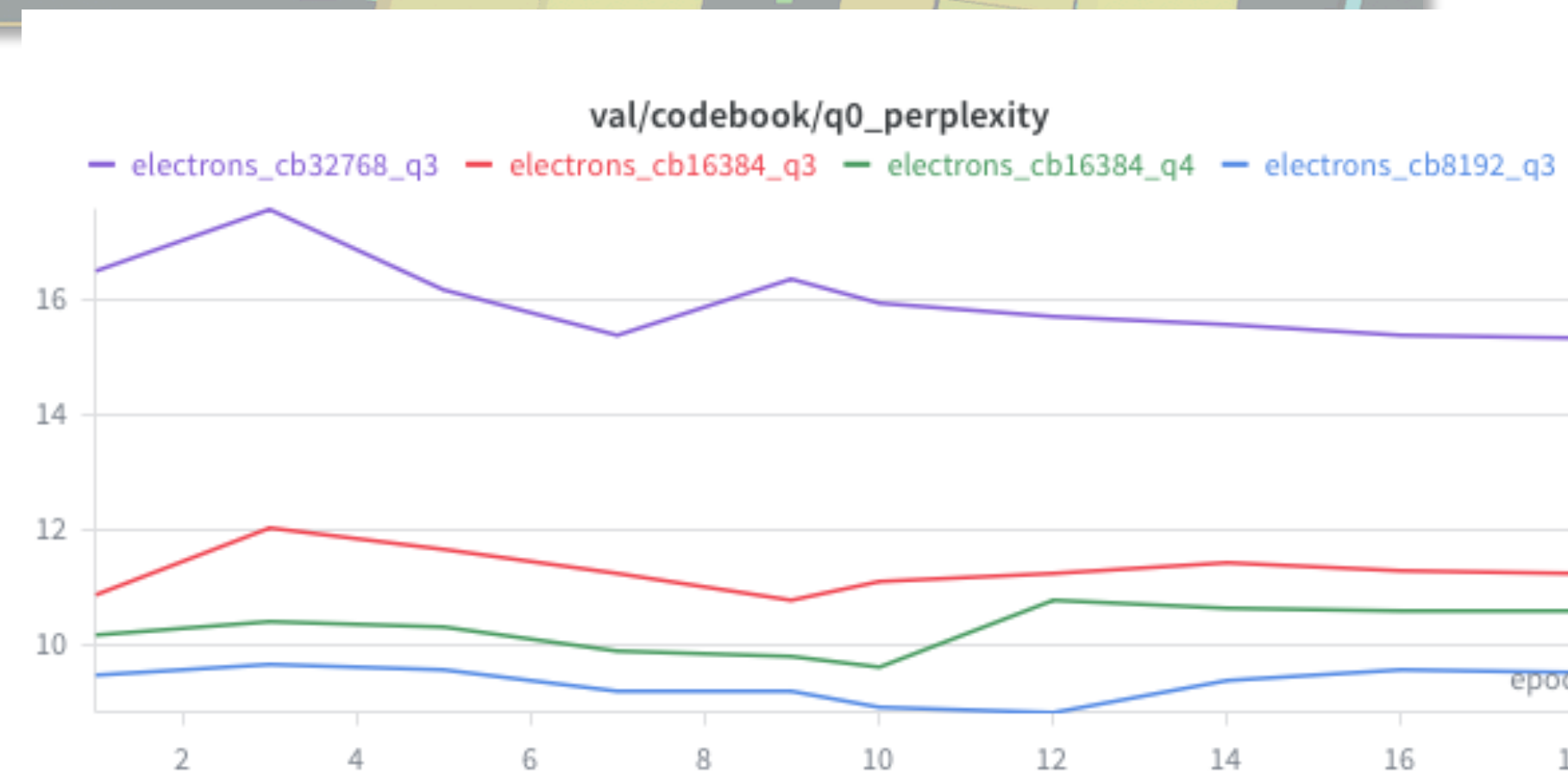
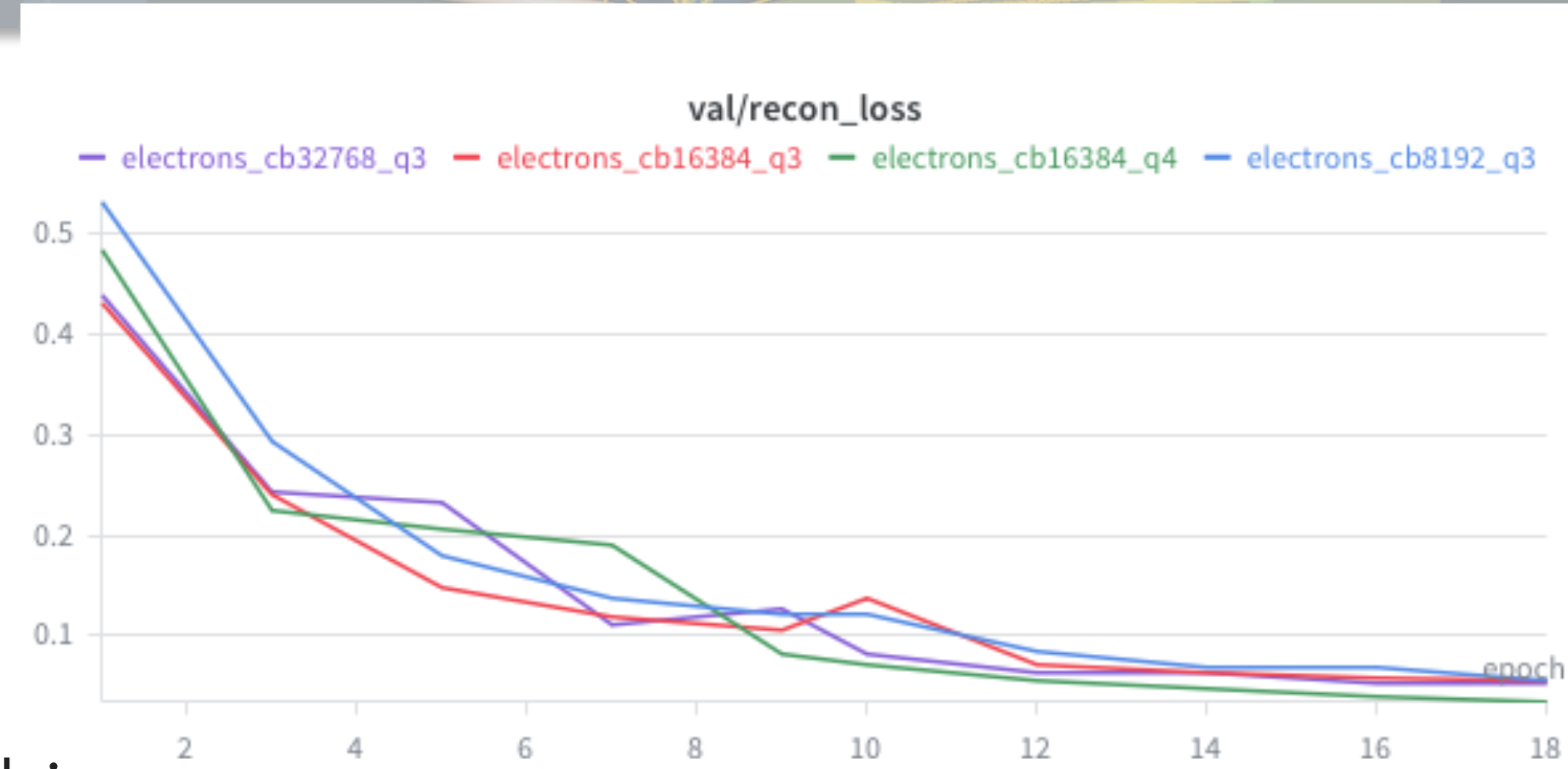
- * HH->4l, ttbar and DY background
- * Data15
- * Event and object info

```
jets      events=8,778,484 objects=16,543,249
electrons events=8,778,484 objects=66,818
muons     events=8,778,484 objects=124,045
photons   events=8,778,484 objects=1,016,327
taus      events=8,778,484 objects=1,678,534
```

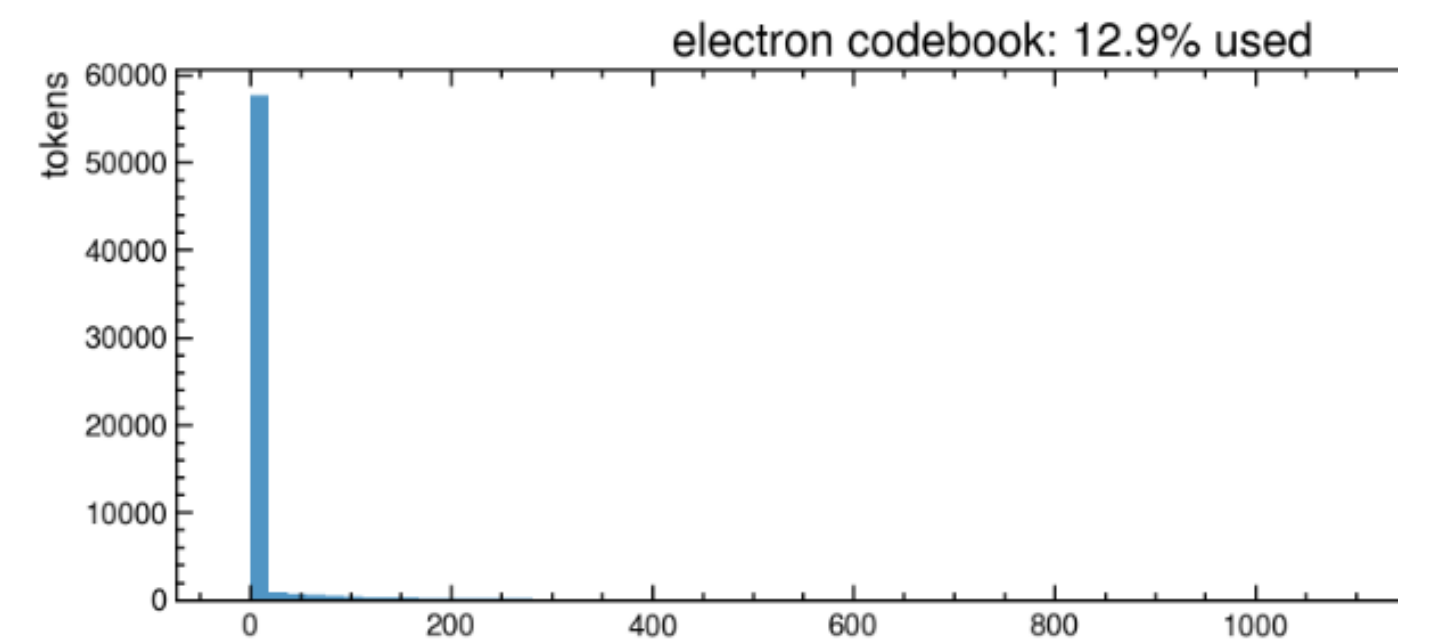
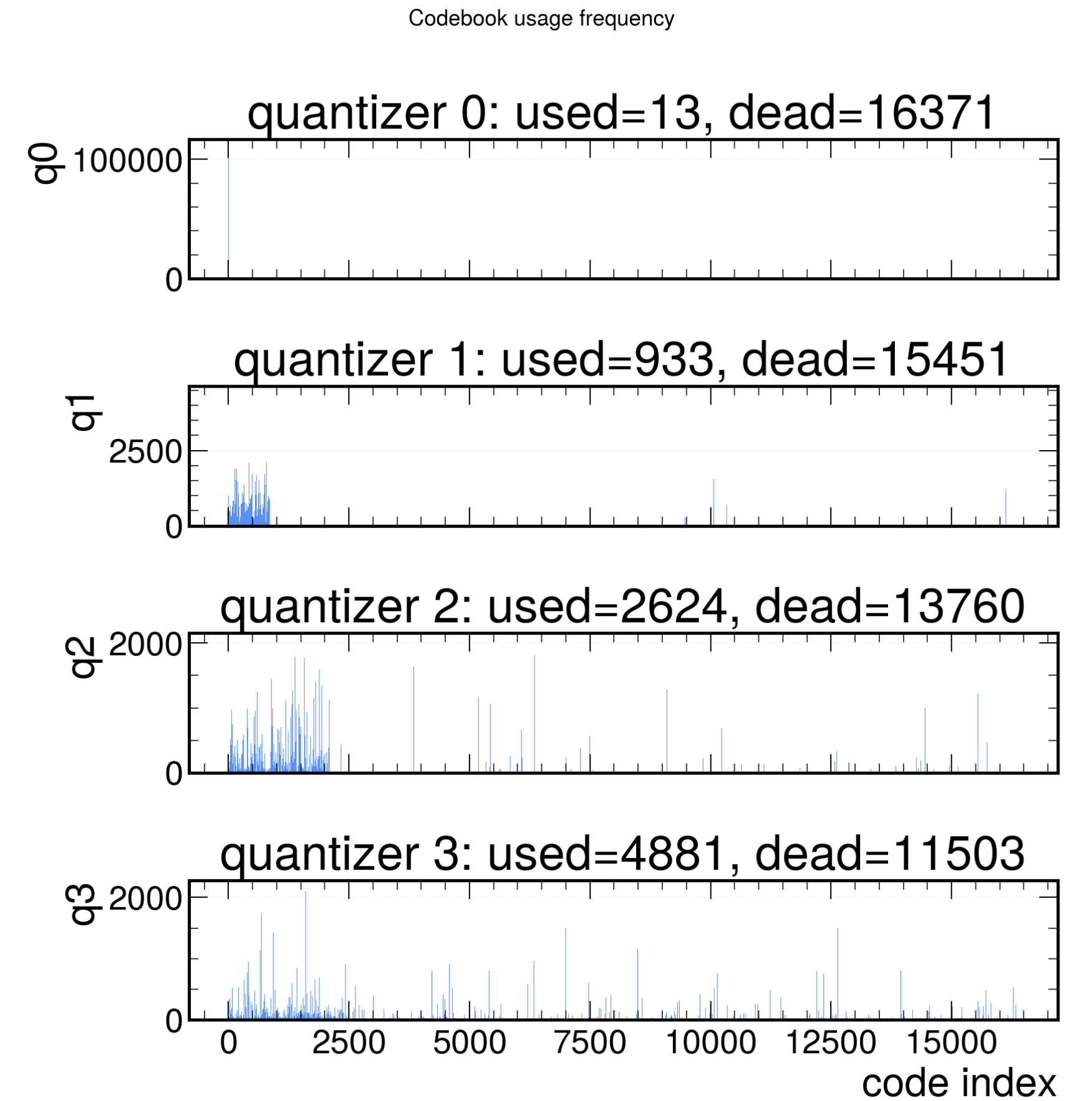
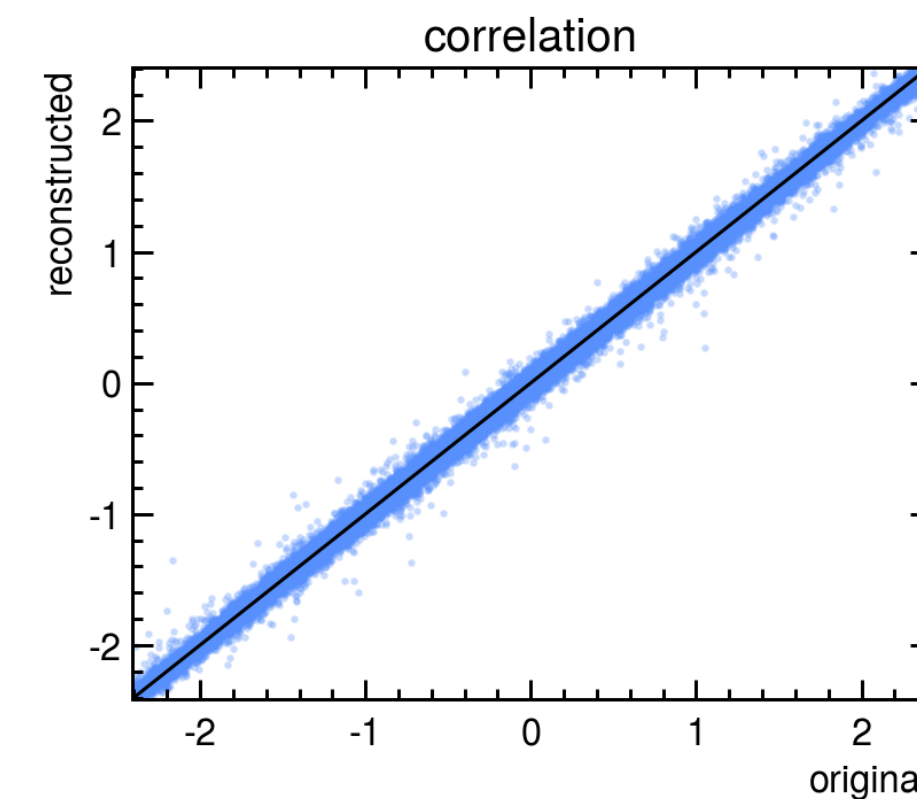
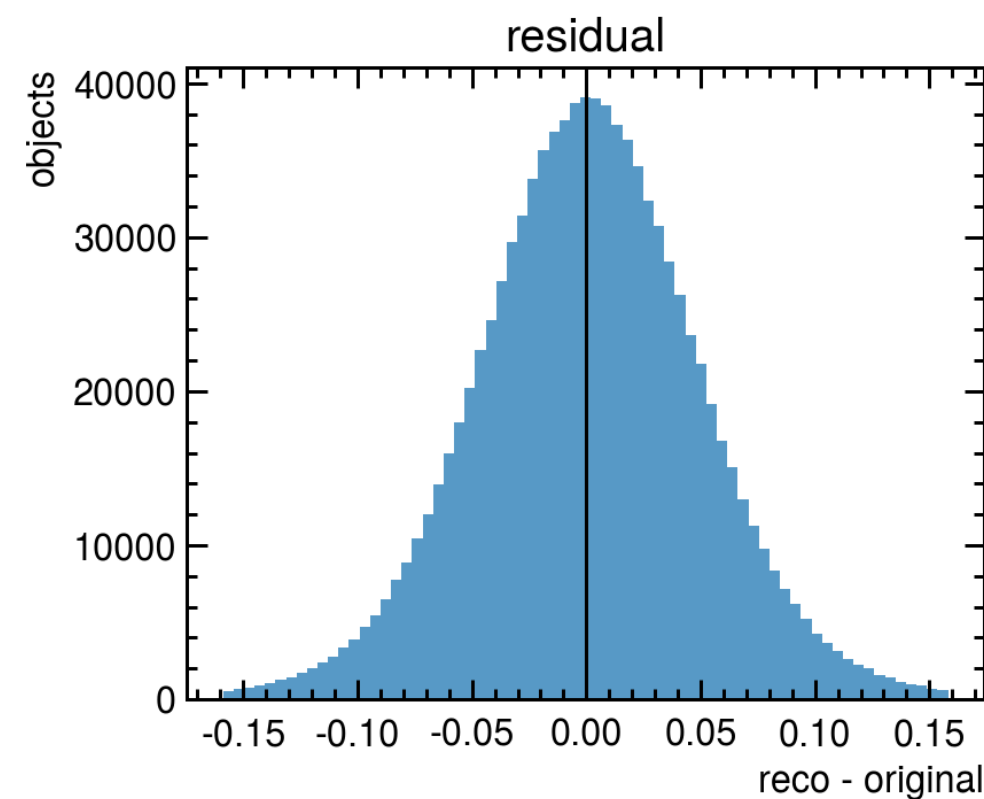
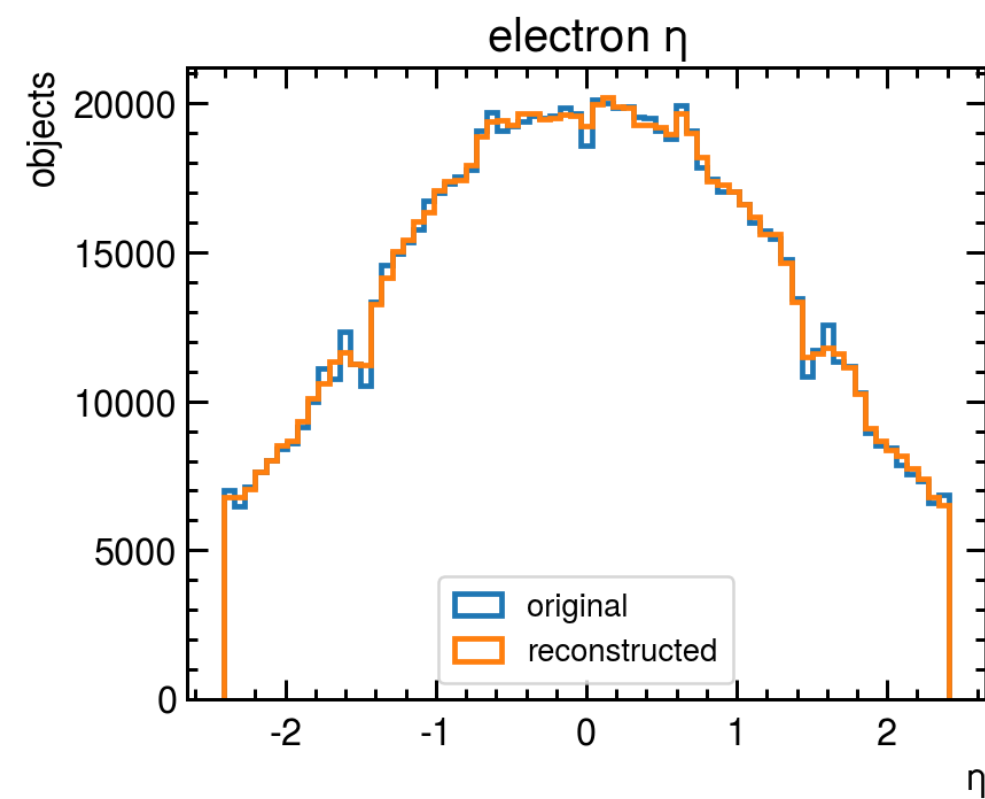
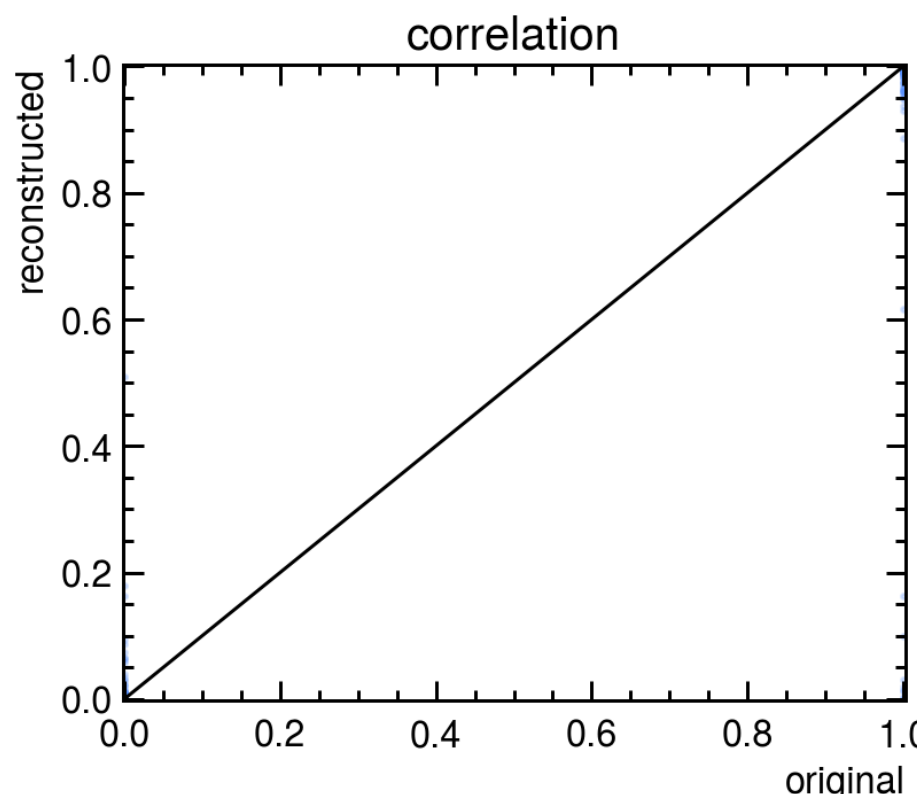
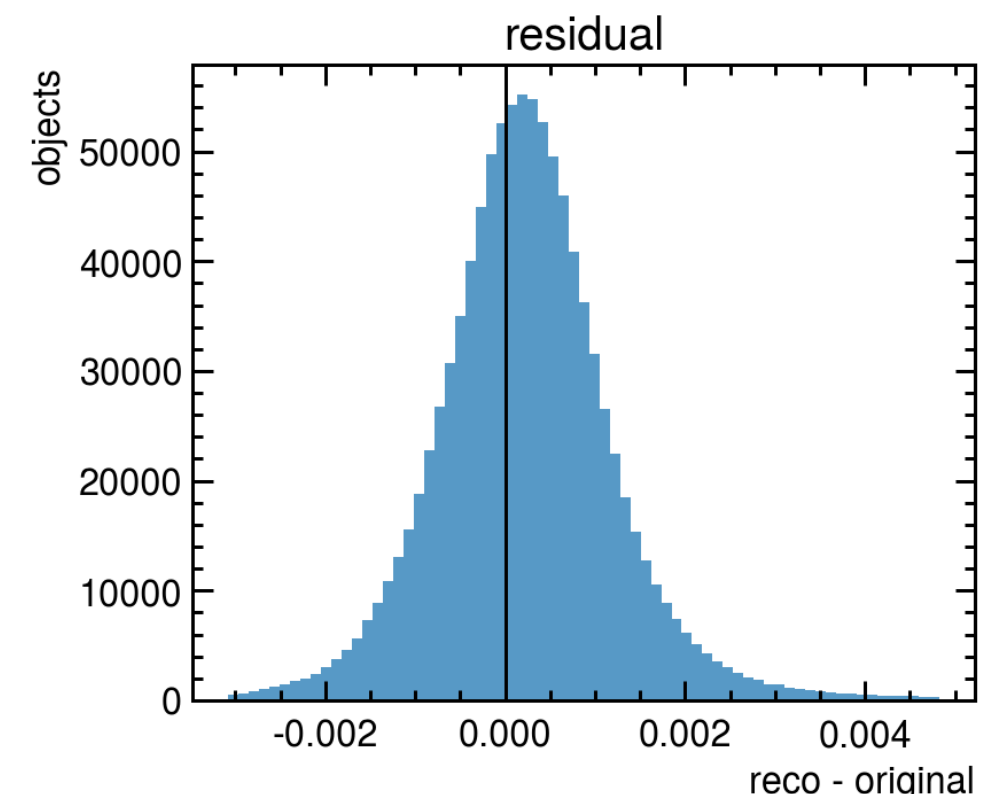
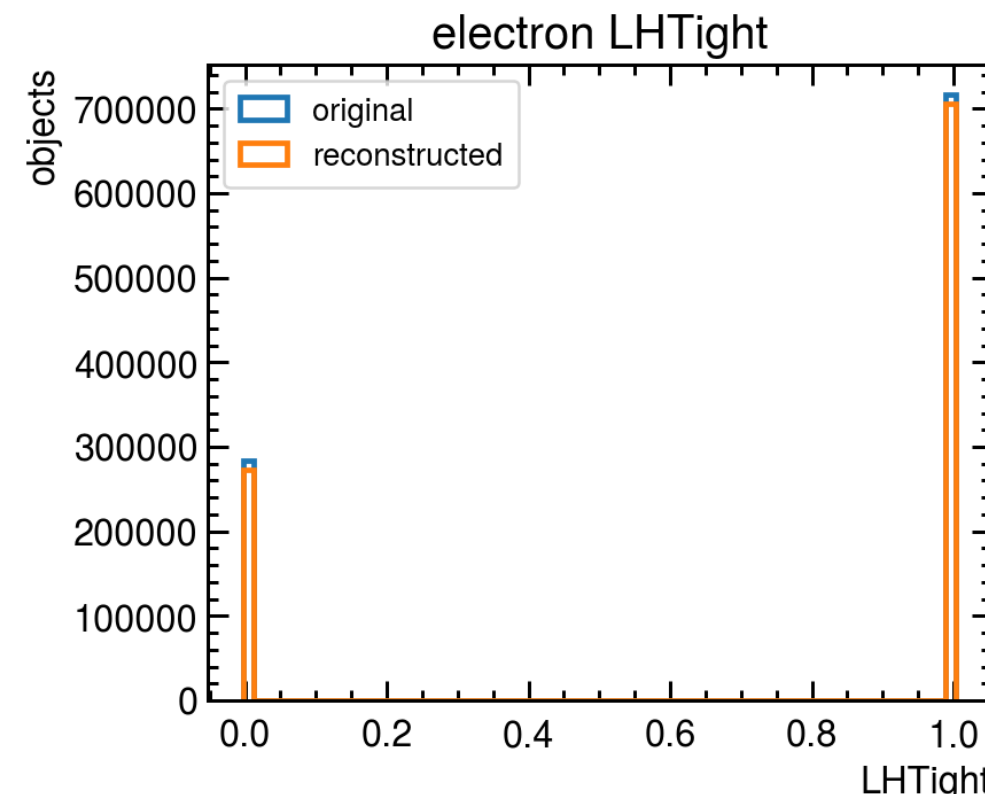
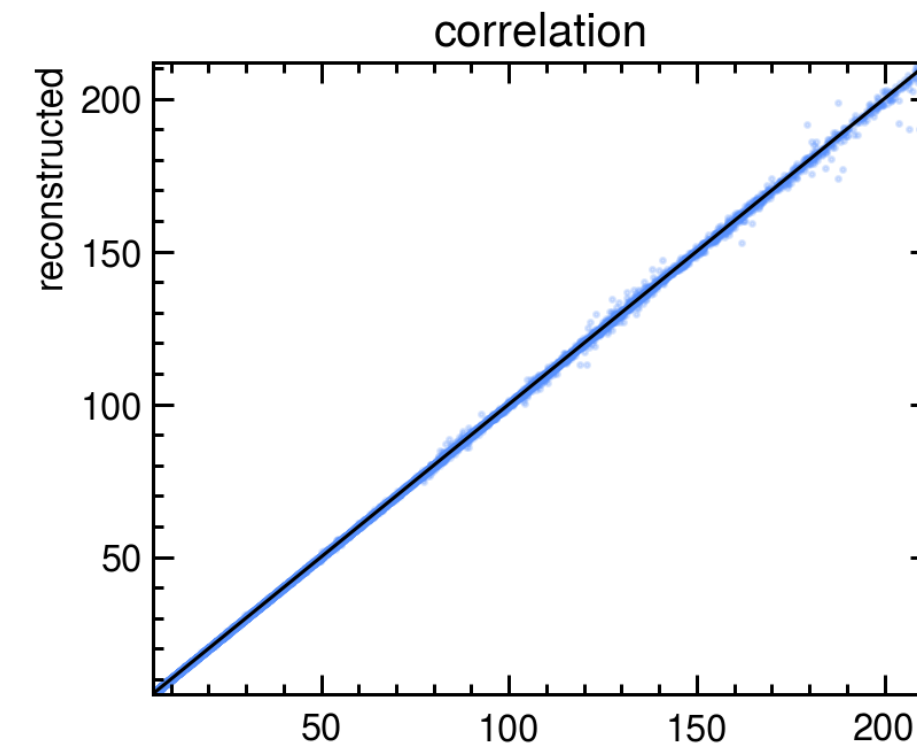
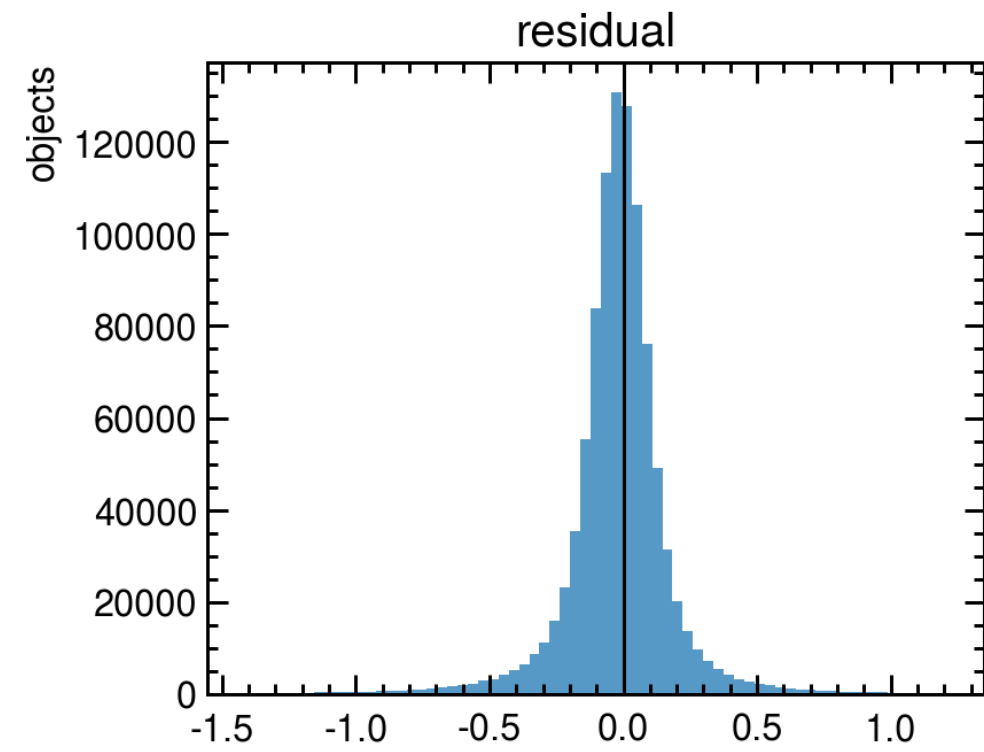
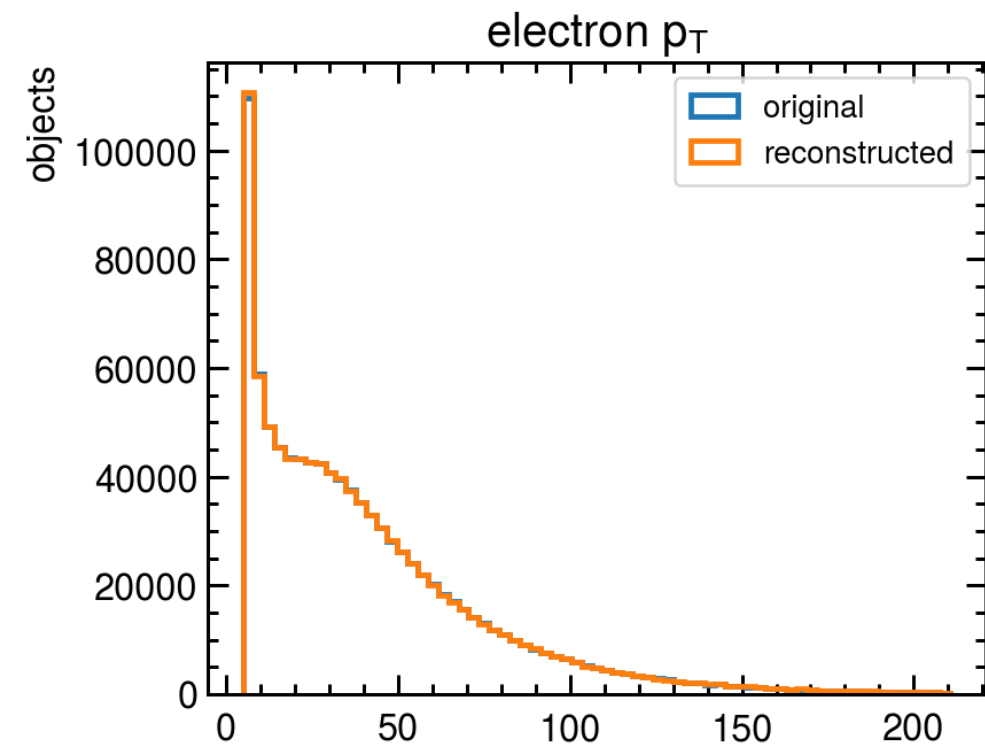
- * Code: <https://github.com/Treasure-AmSC/heptokens/tree/event-object-tokenization>

Electrons

- * For electrons, the result is pretty clear: **electrons_cb16384_q4** is best overall.
- * It has the lowest or near-lowest val/recon_loss.
- * It clearly improves pt, eta, phi, and charge.
- * It improves the ID variables, especially LHMedium and LHTight.
- * The 4th quantizer q3 is active, with ~650-680 used codes, so it is actually contributing.

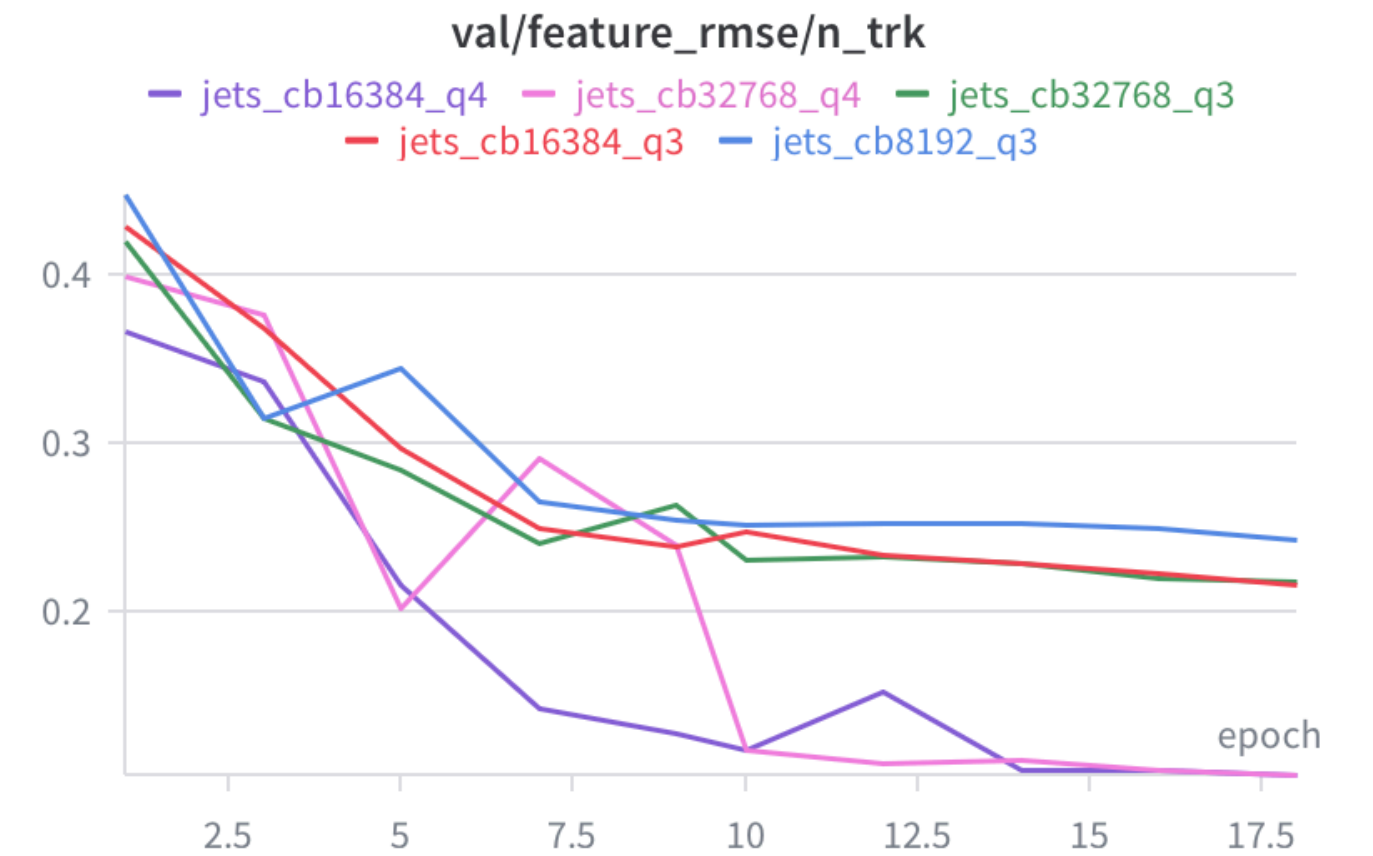
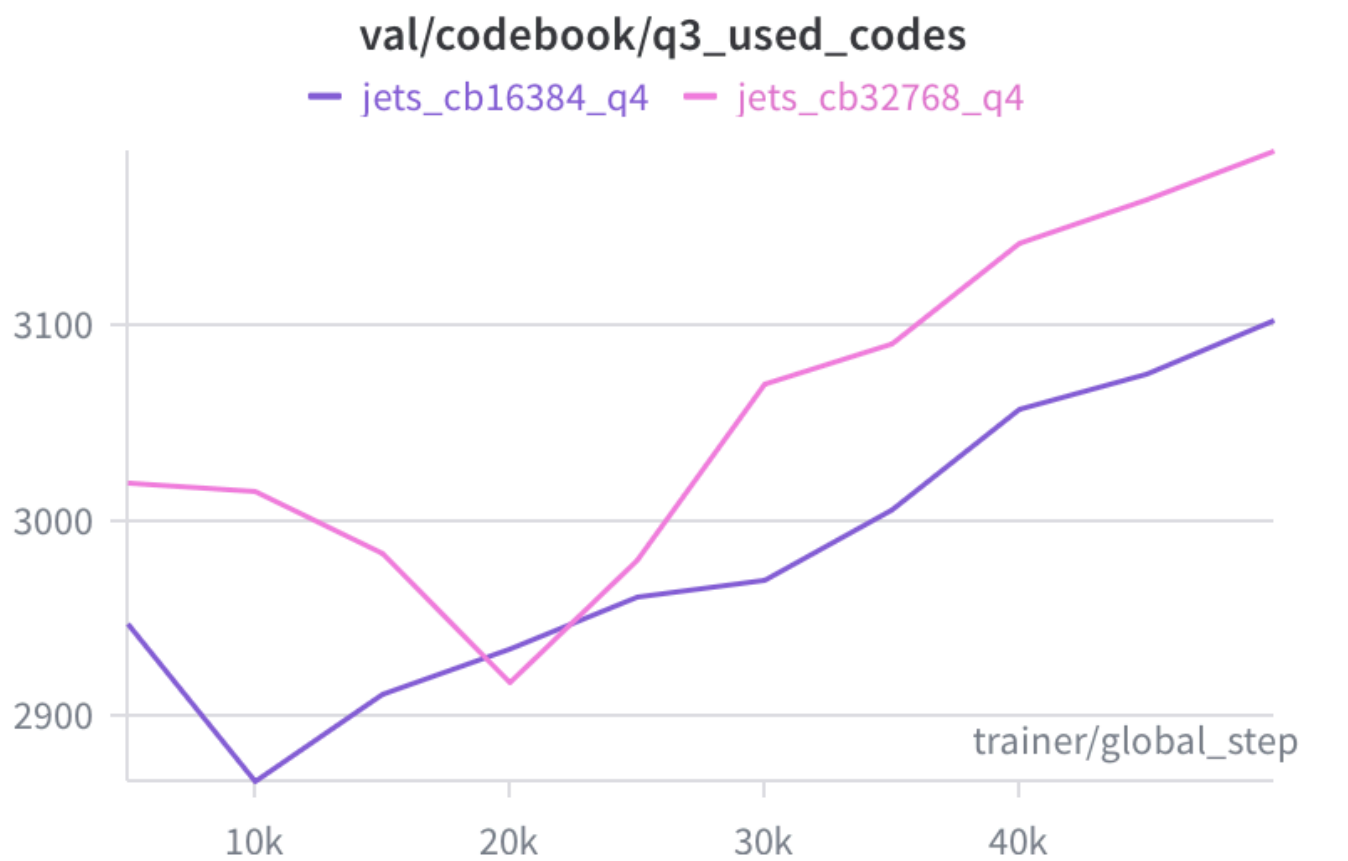
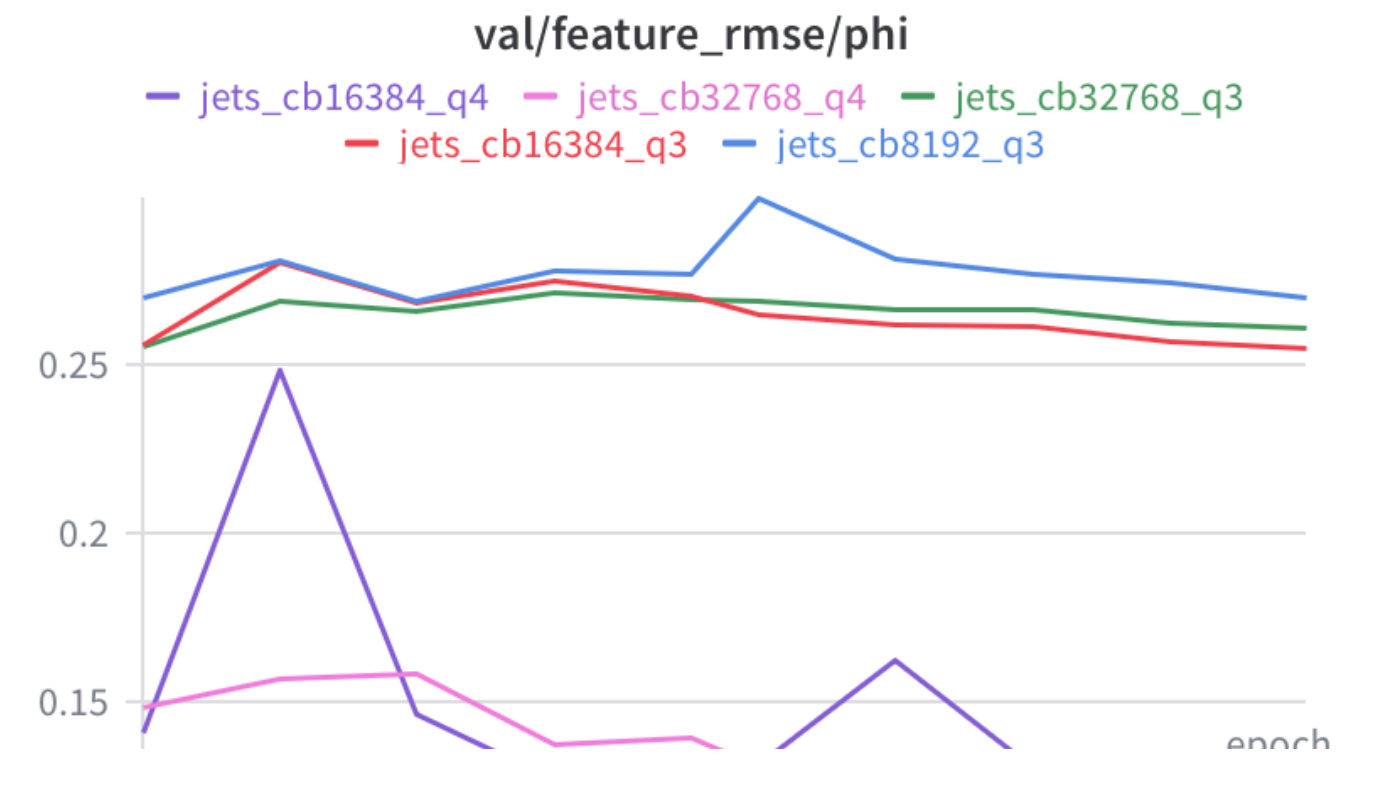
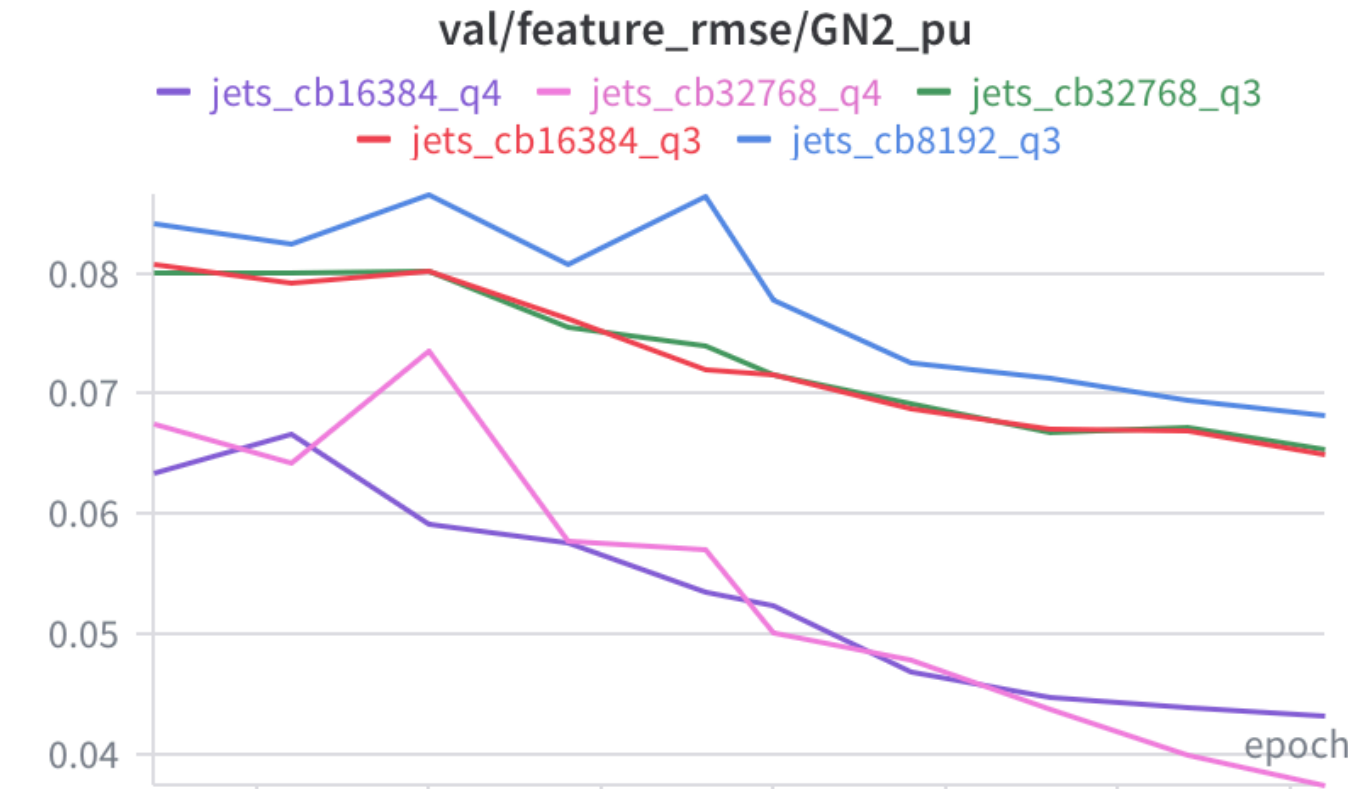
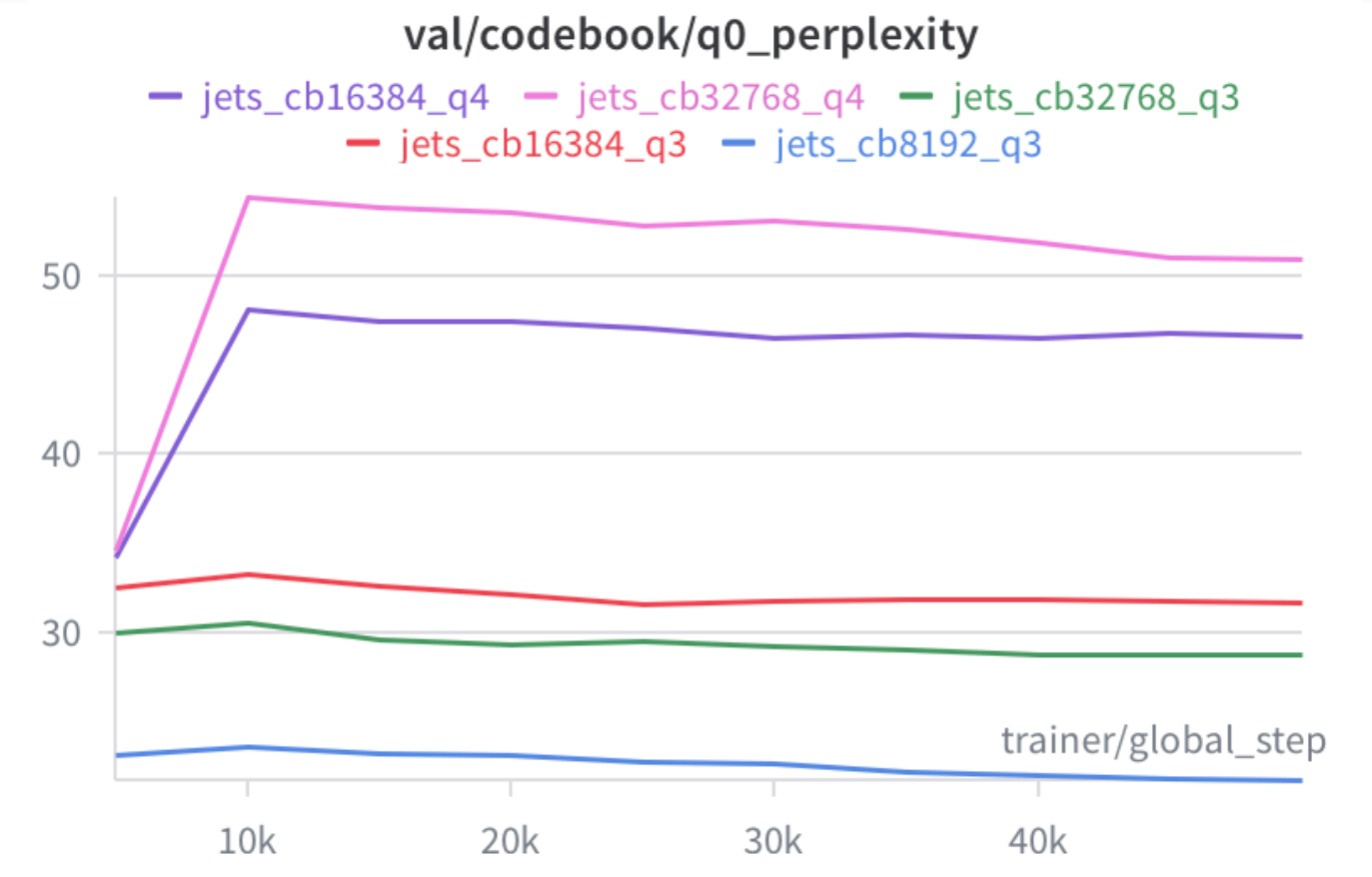
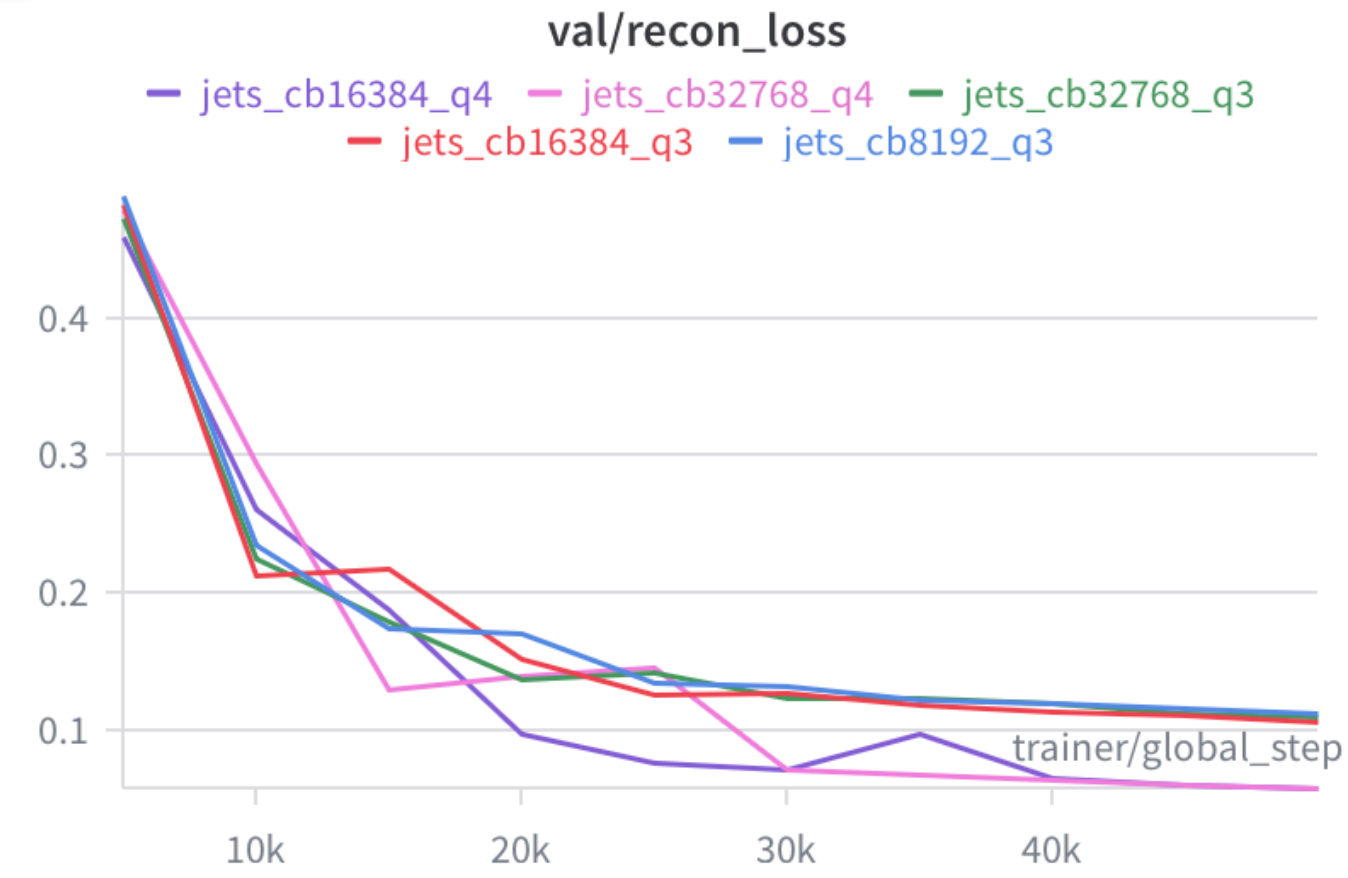


Reco - electrons

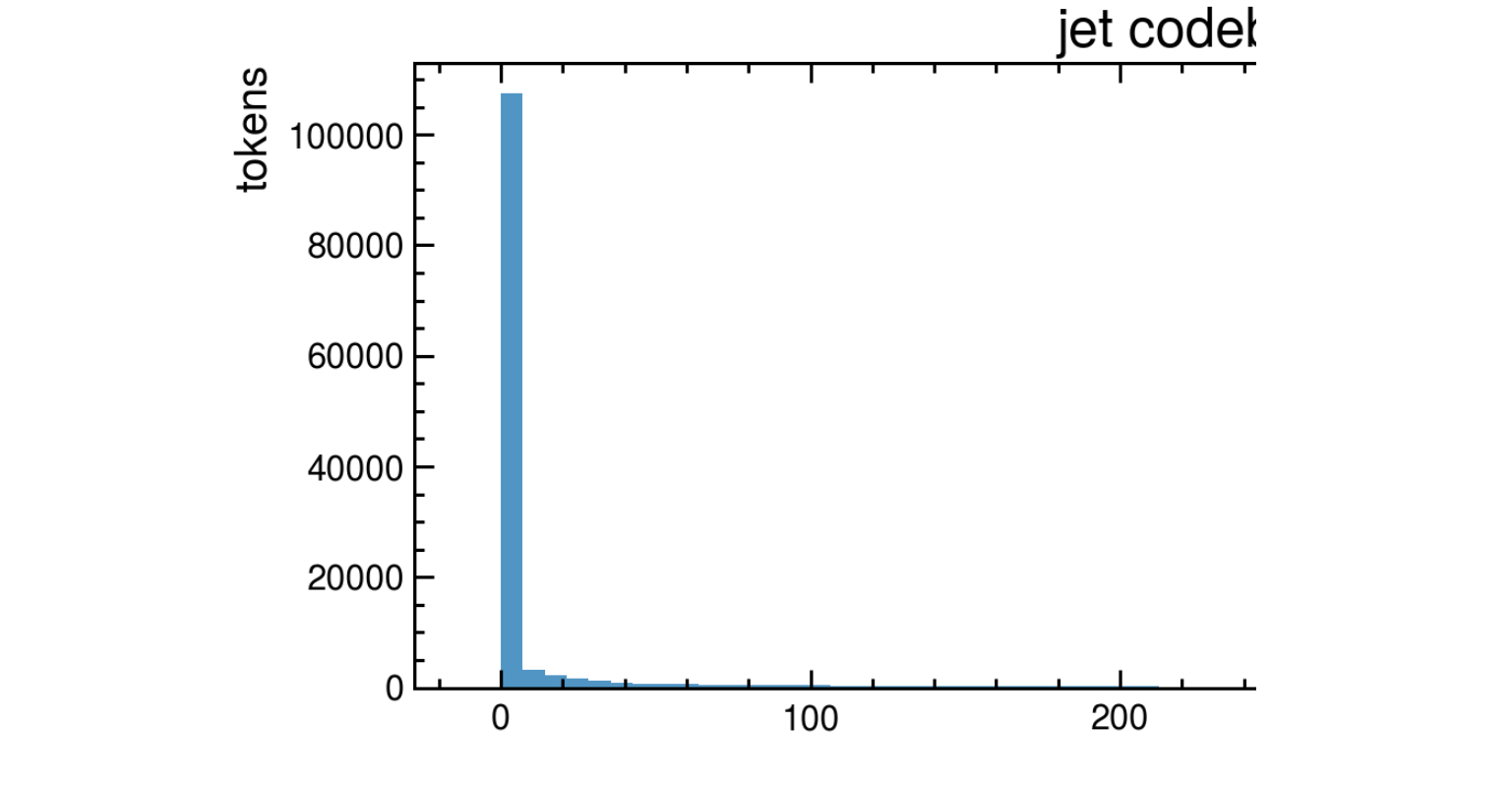
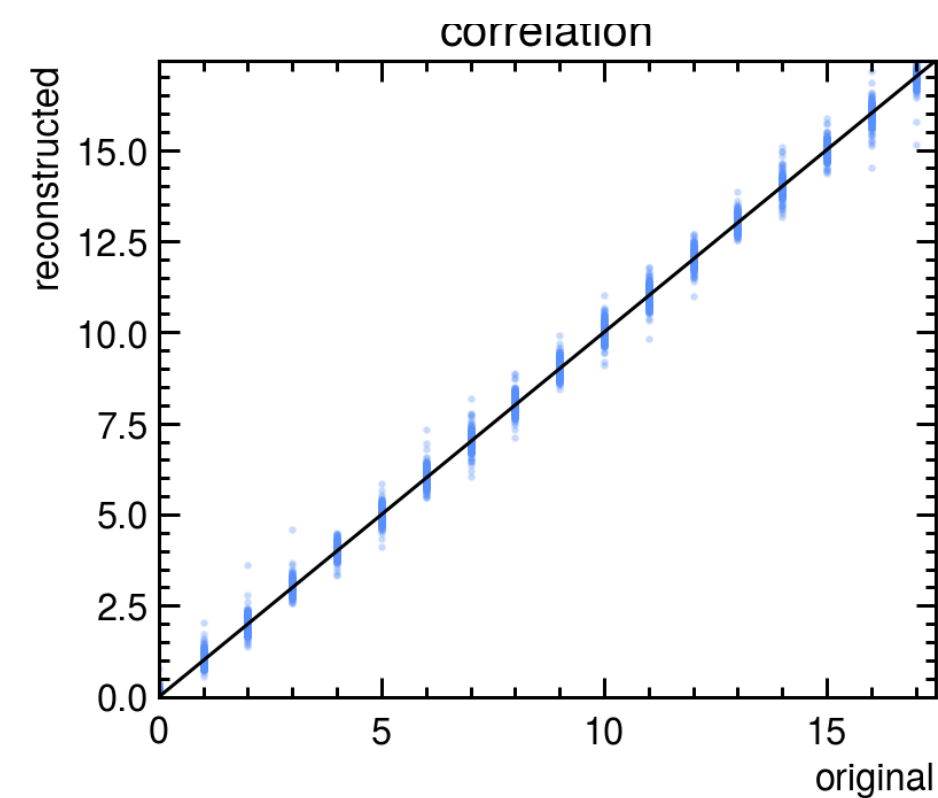
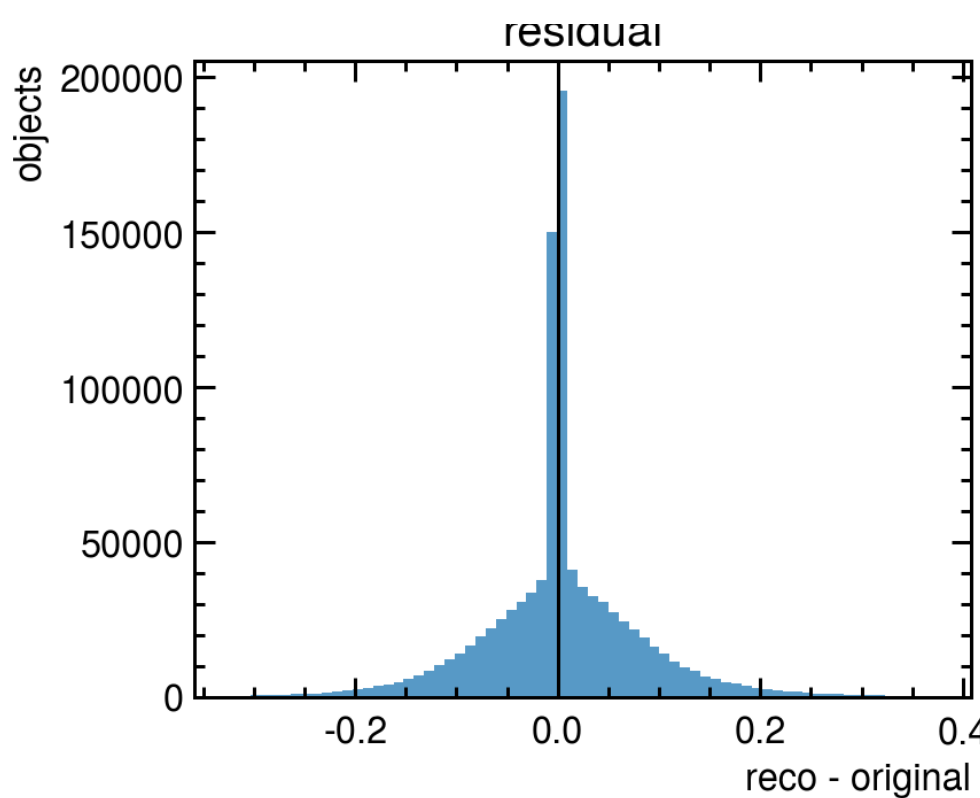
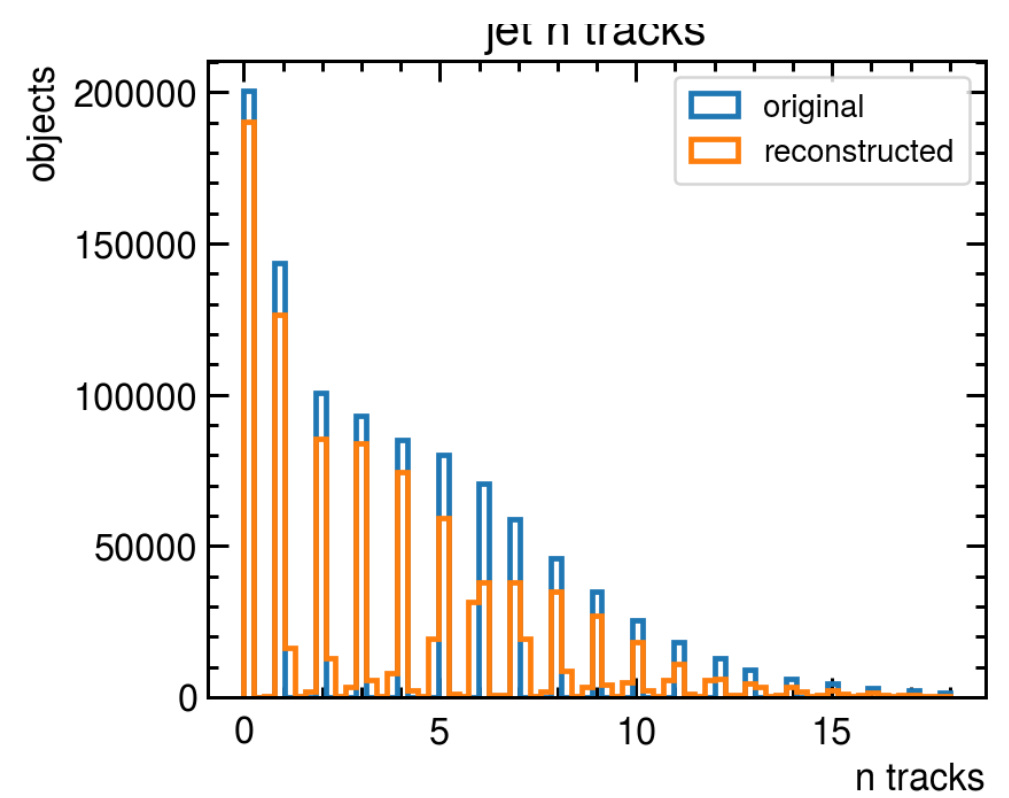
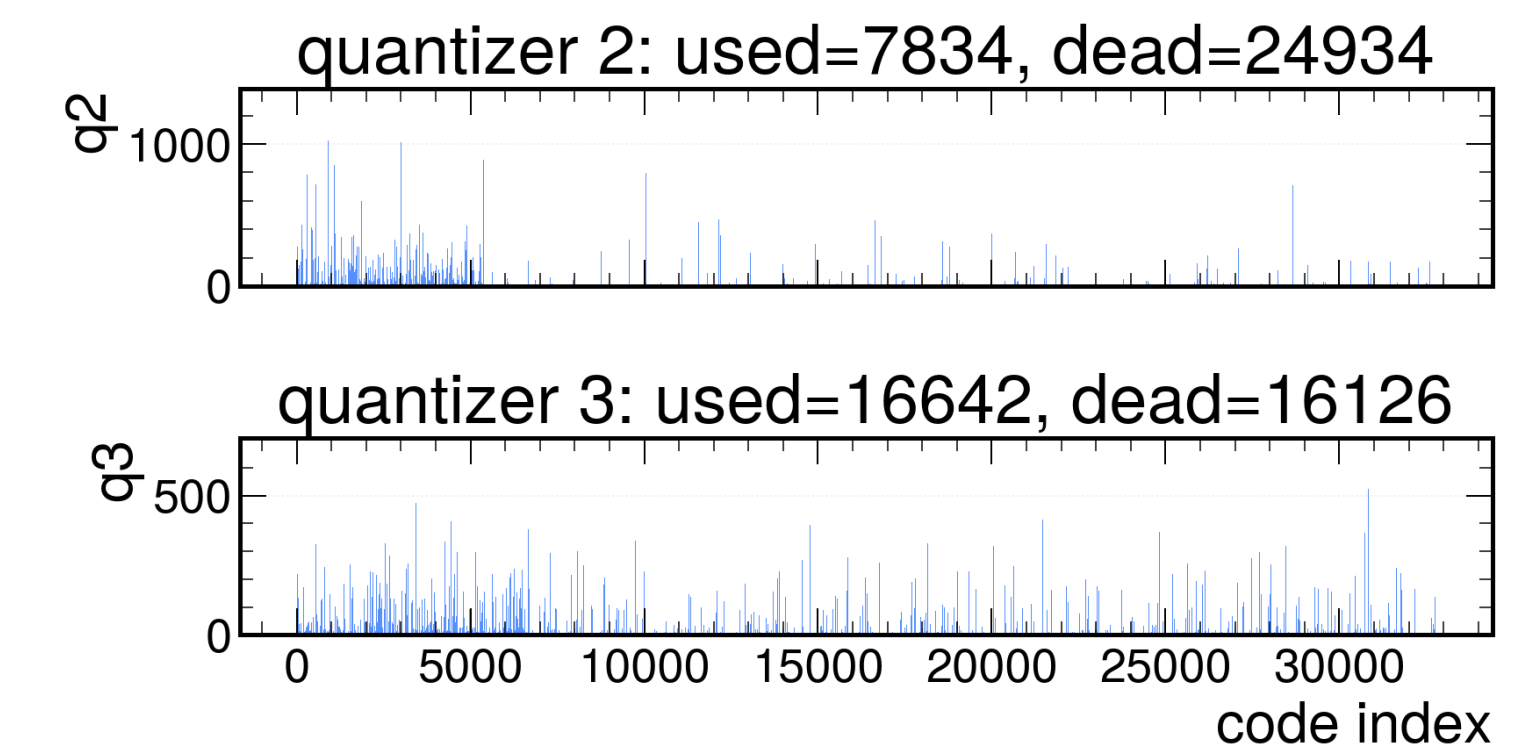
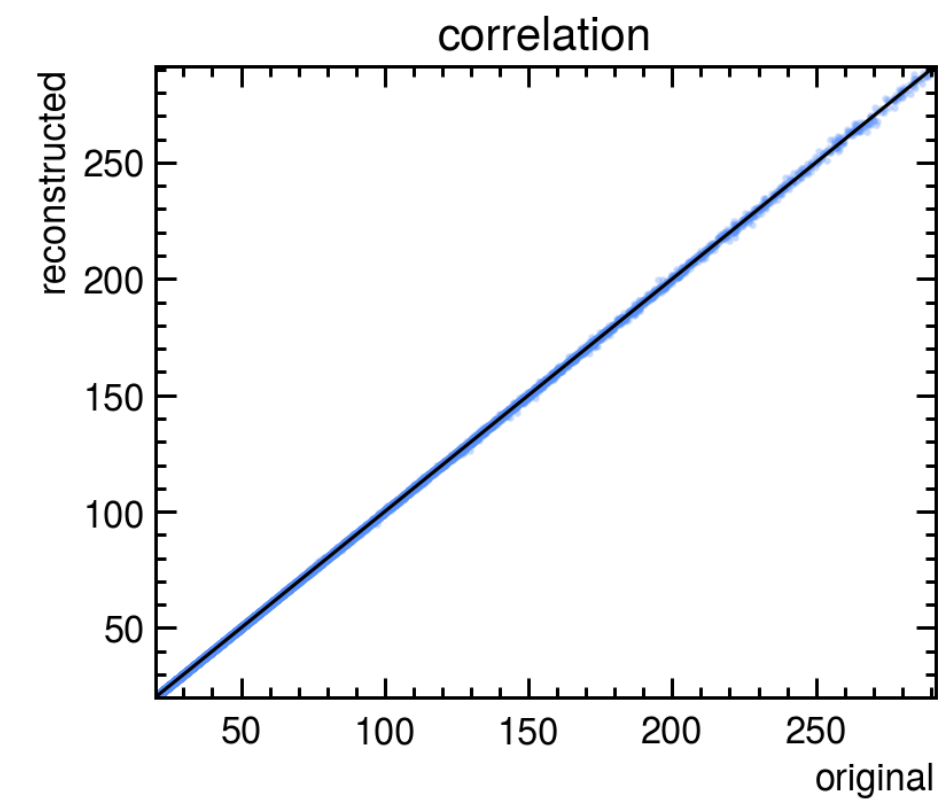
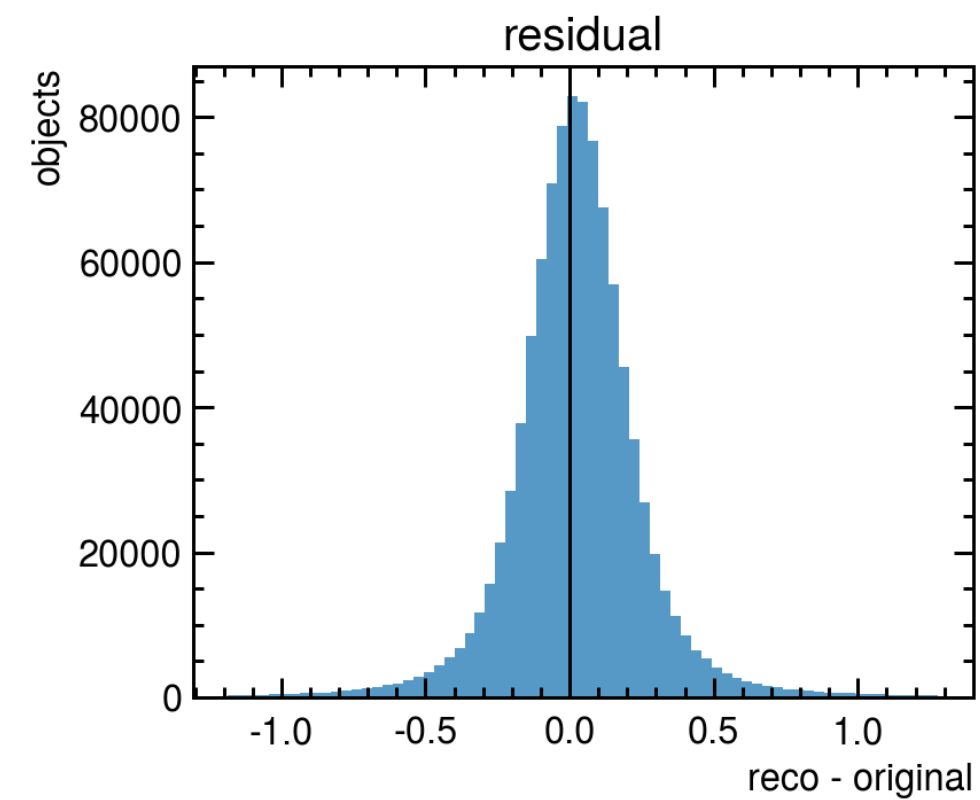
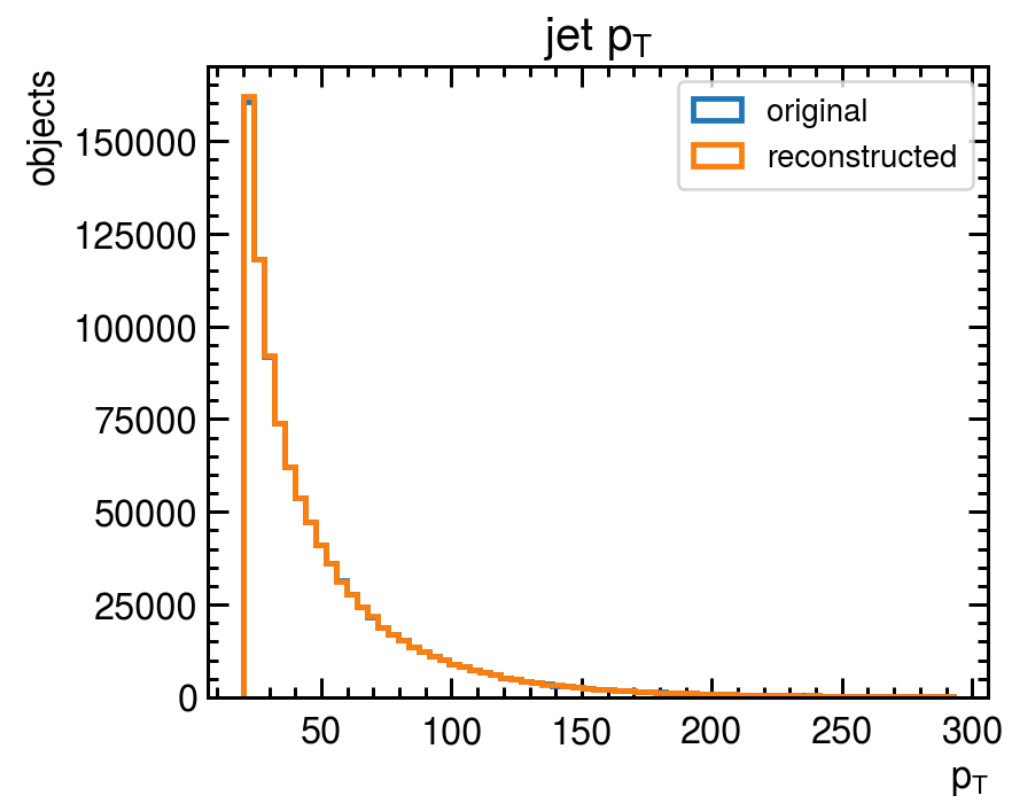
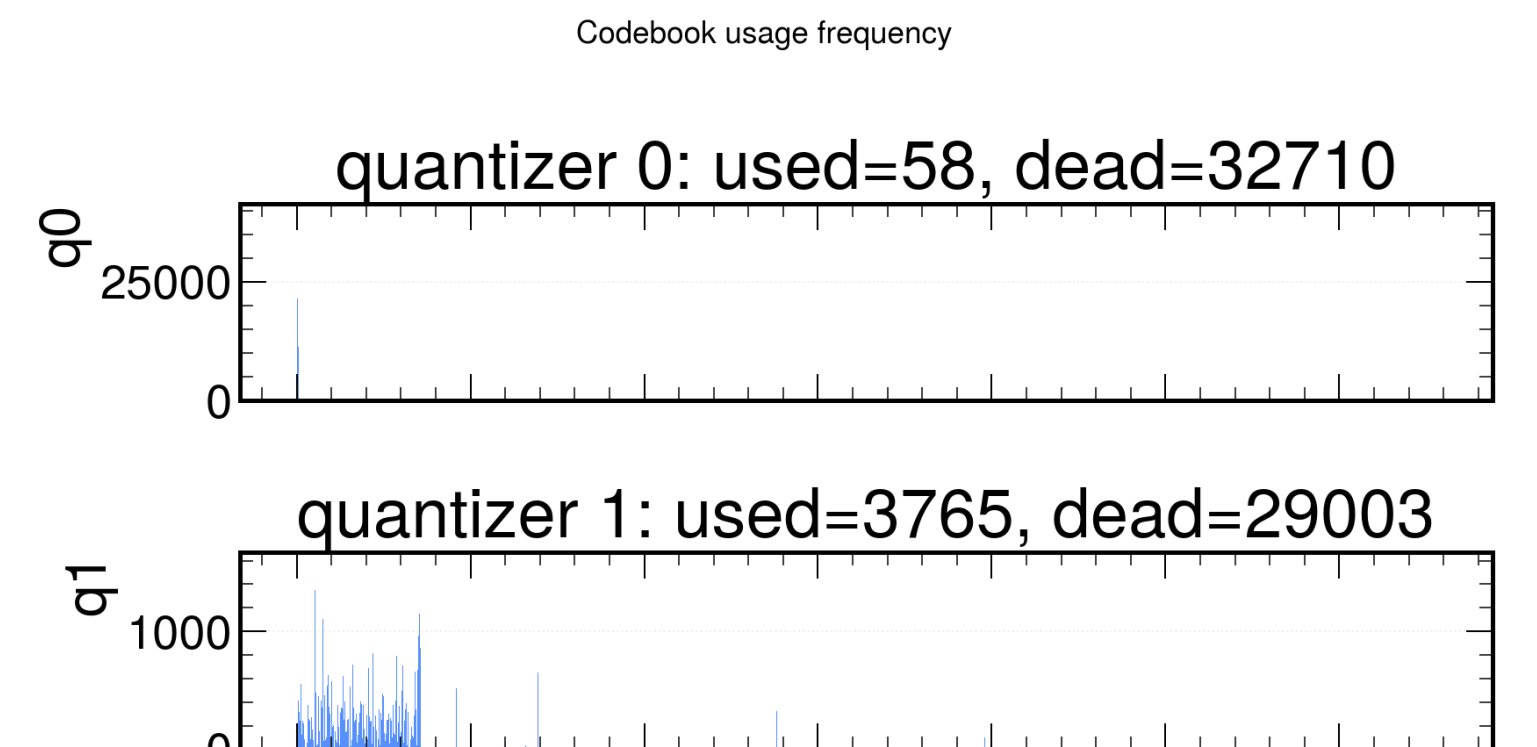
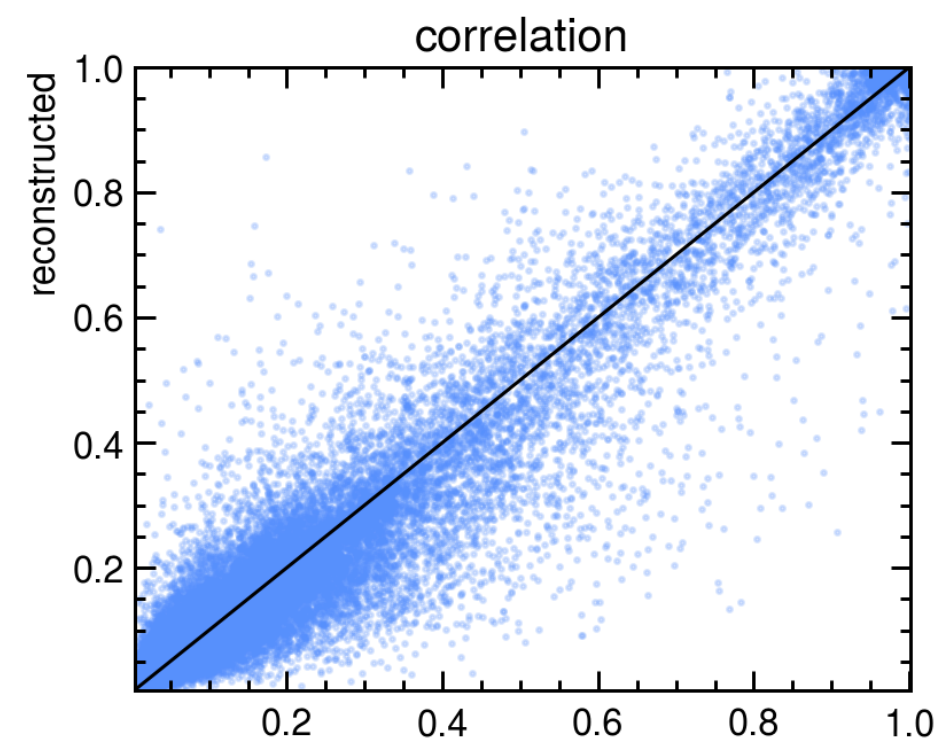
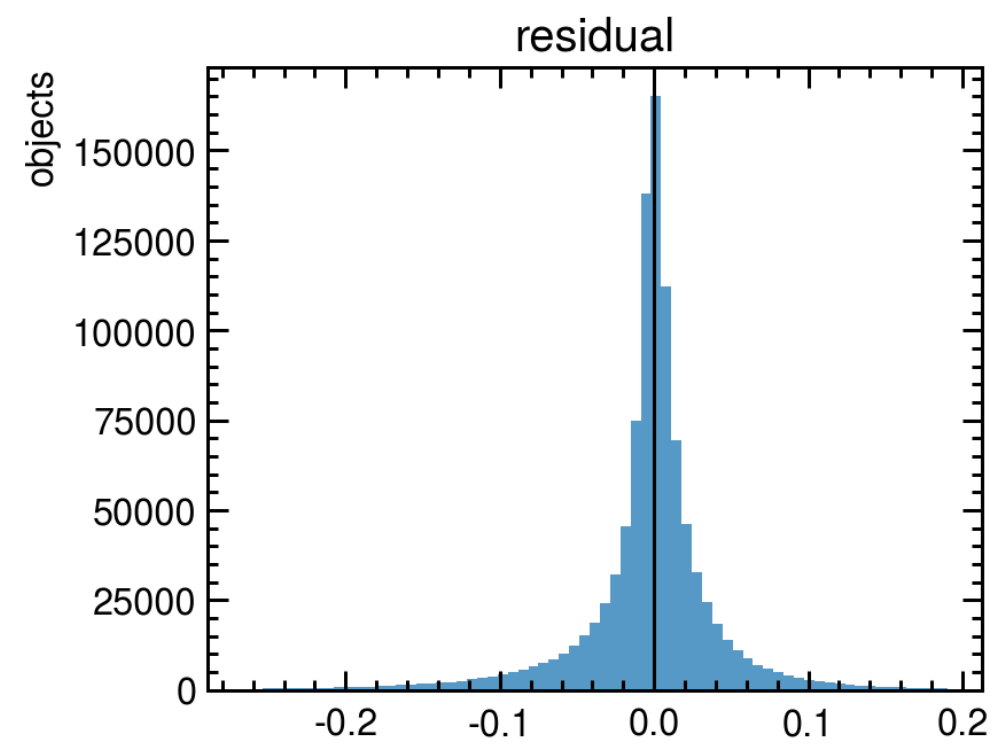
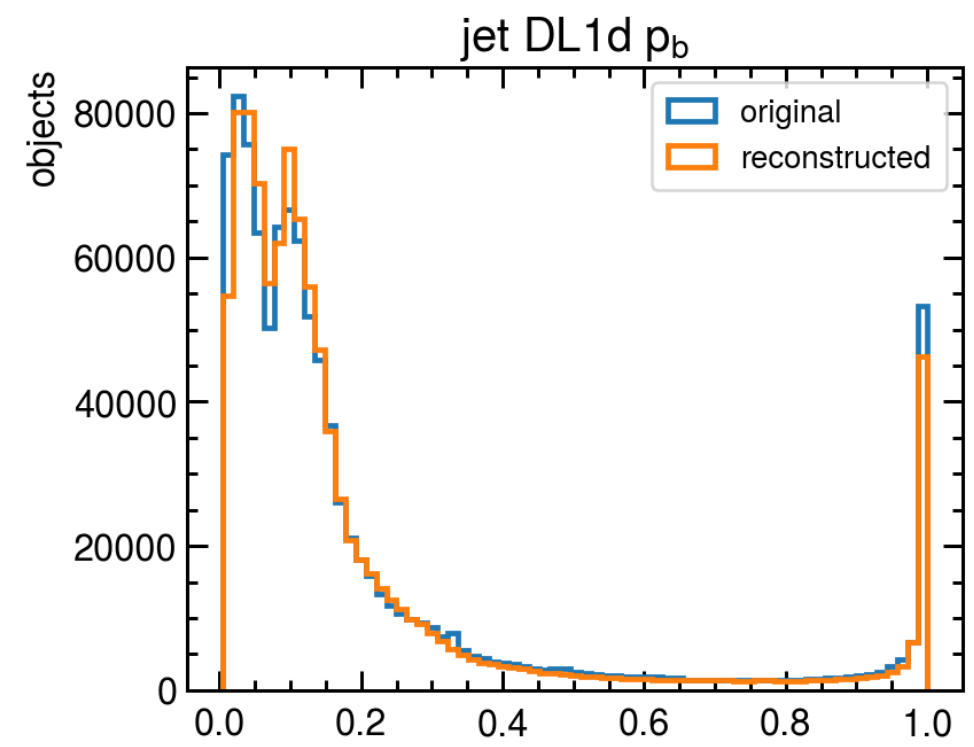


Jets

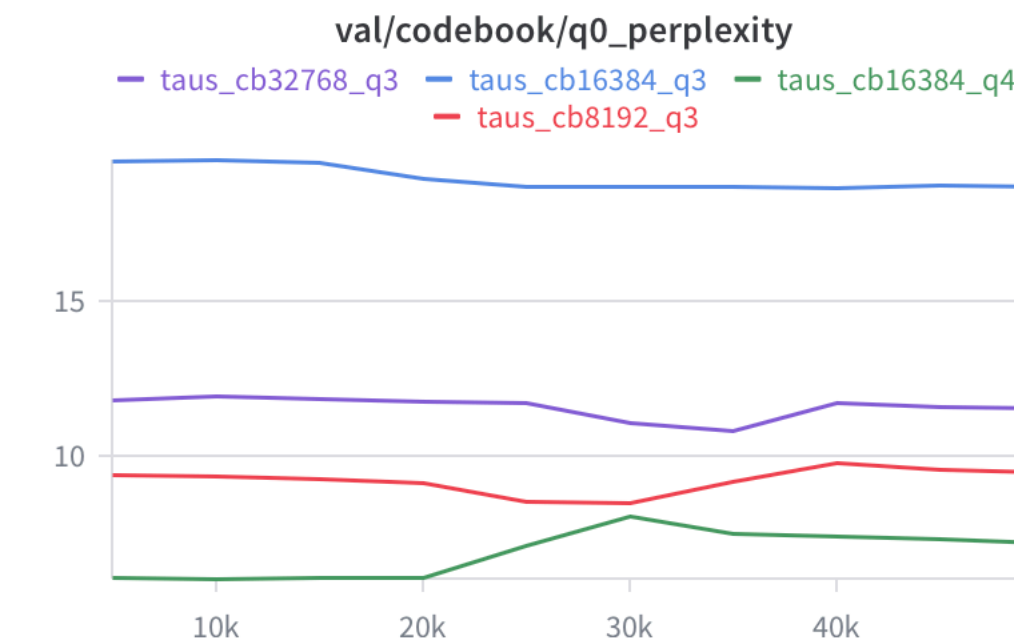
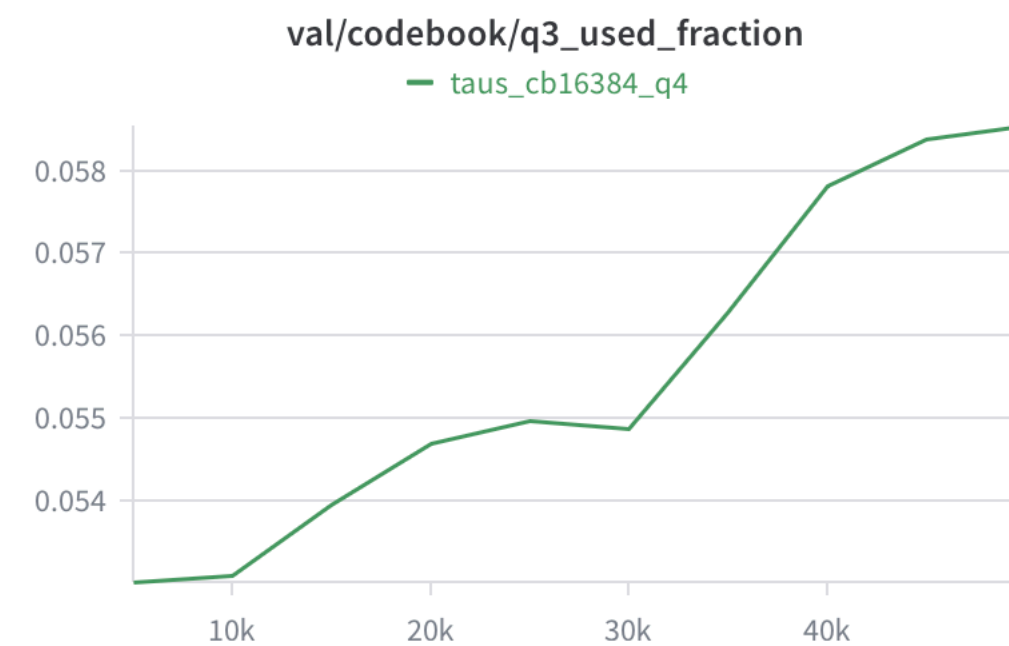
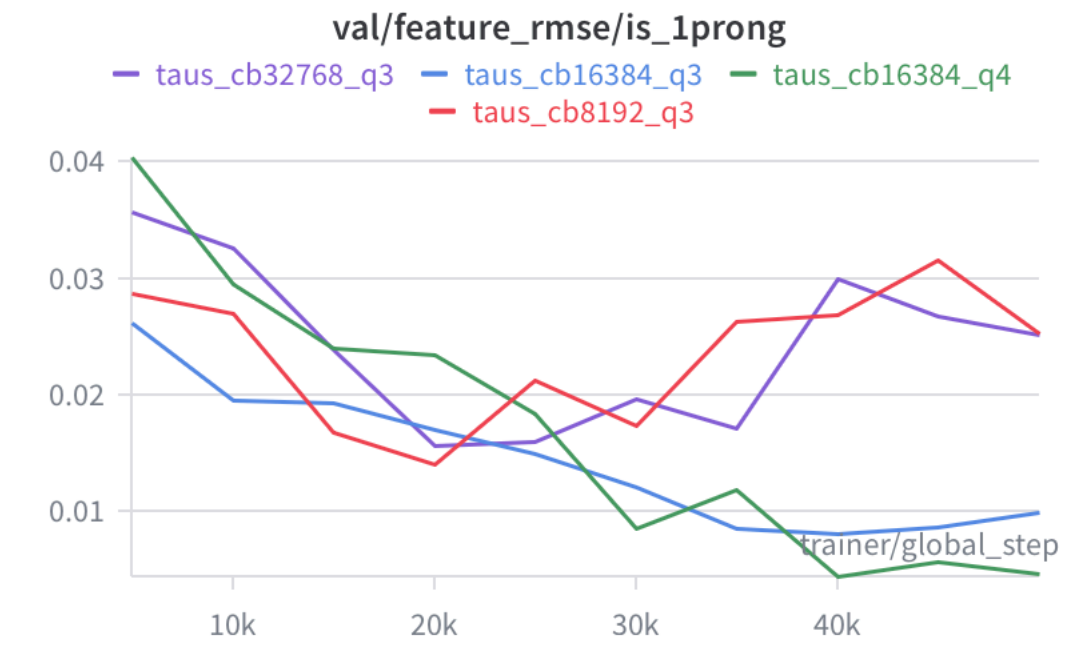
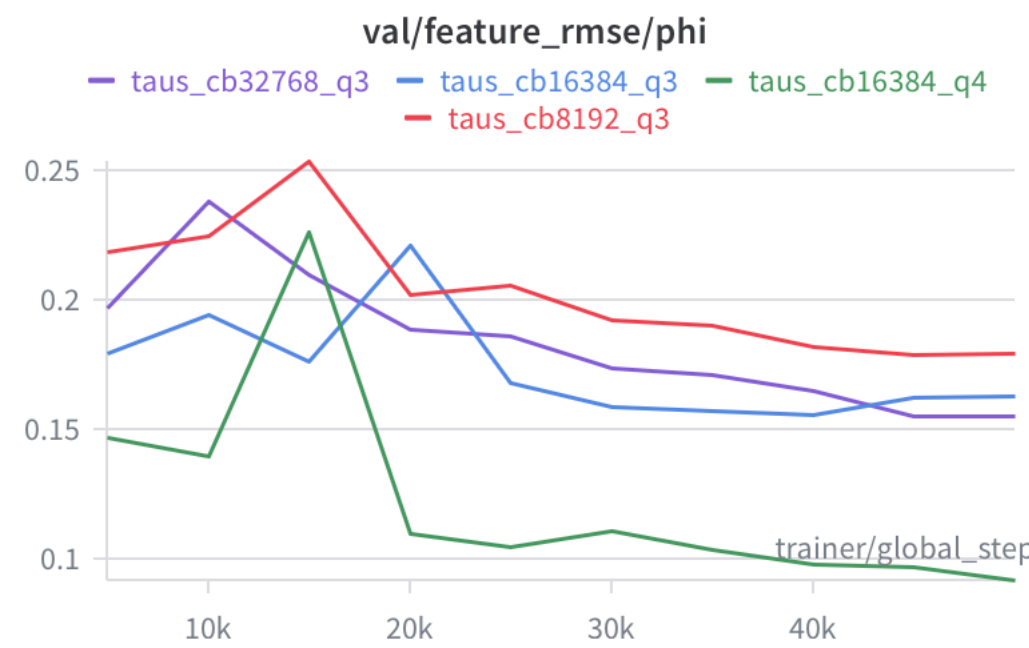
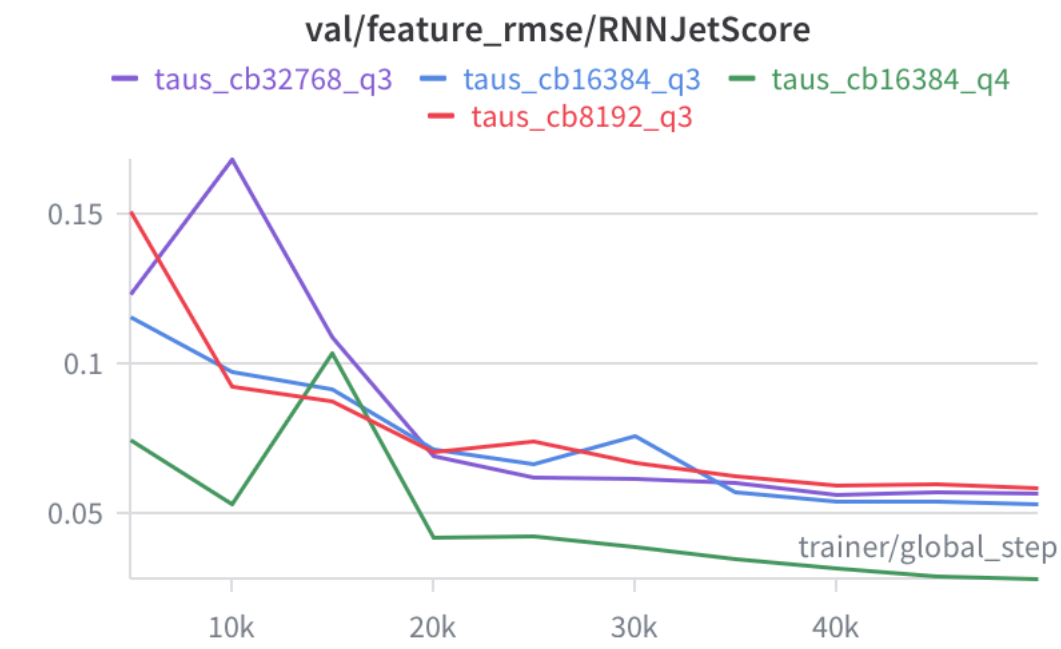
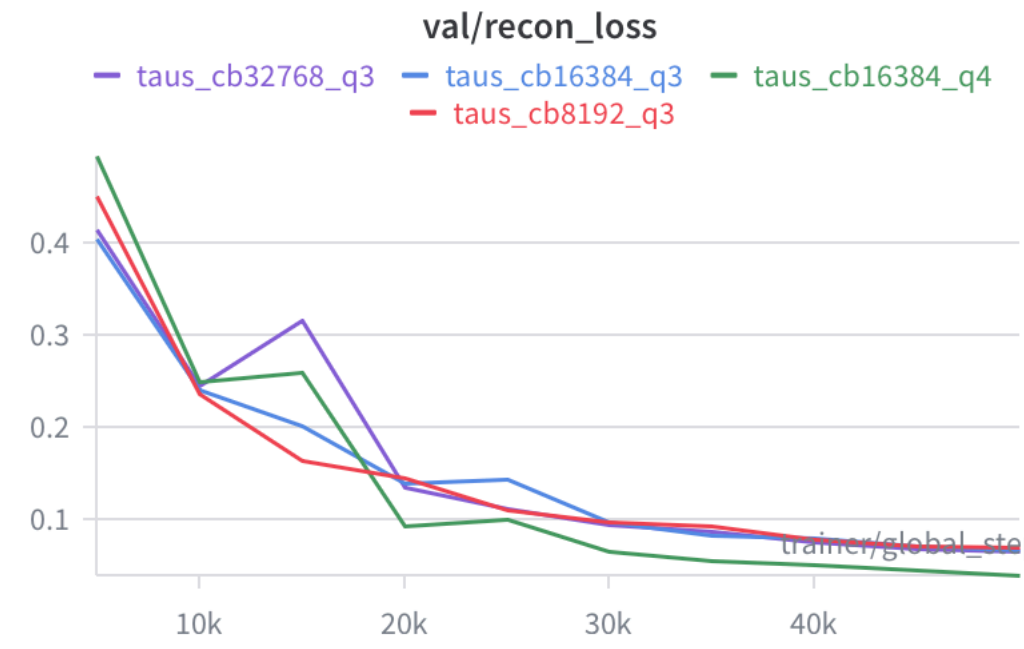
- * 32k q4 is not clearly better
- * It Improves over the q3 runs, but it usually does not beat 16384 q4
 - ▶ For q4, q3 is active and its perplexity increases, so the 4th quantizer is actually helping
- * It also uses a smaller fraction of the codebook
- * 8192 q3 is worse
- * q0 is still collapsed for everything
- * Best reco: **16384_q4**



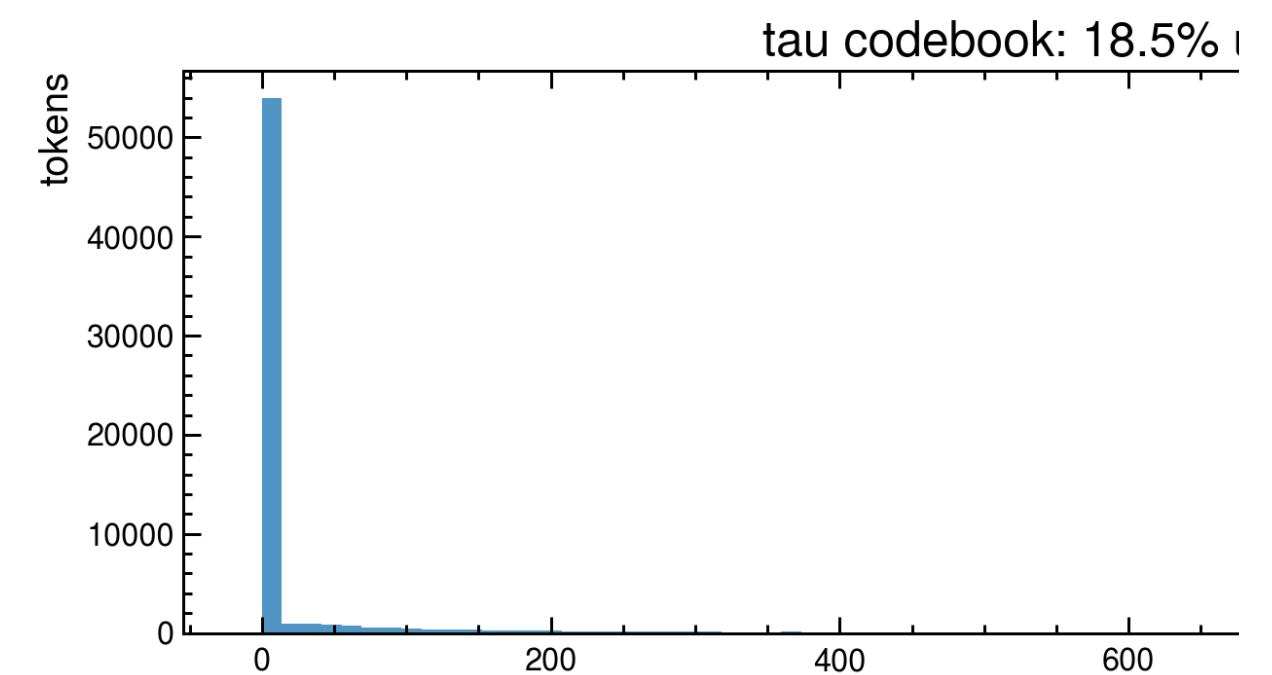
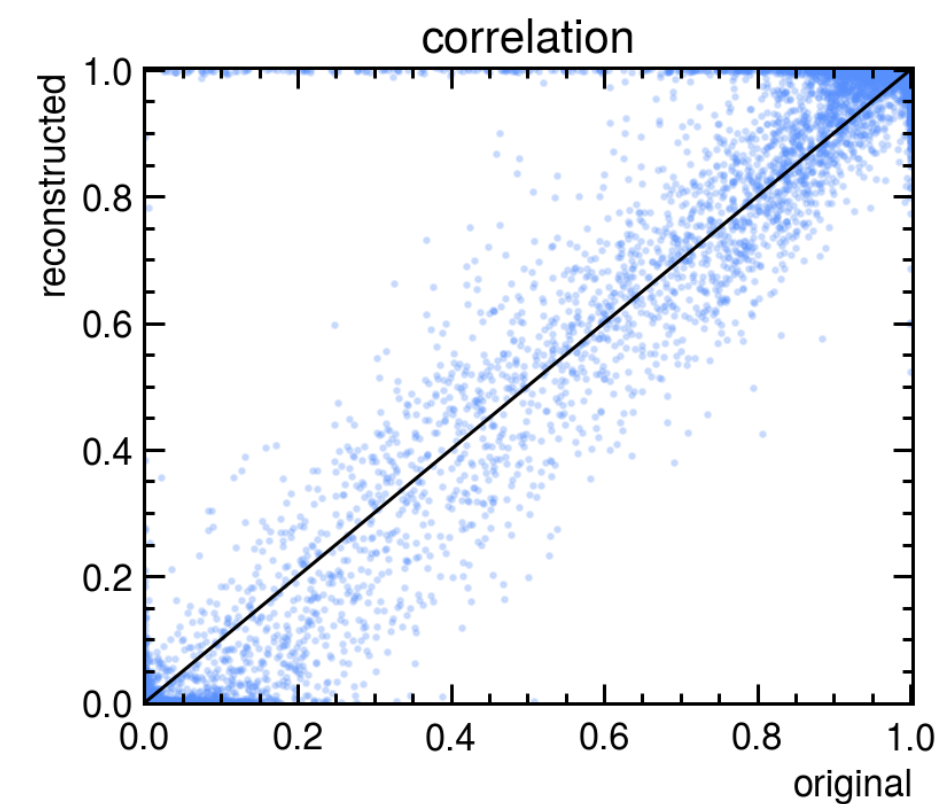
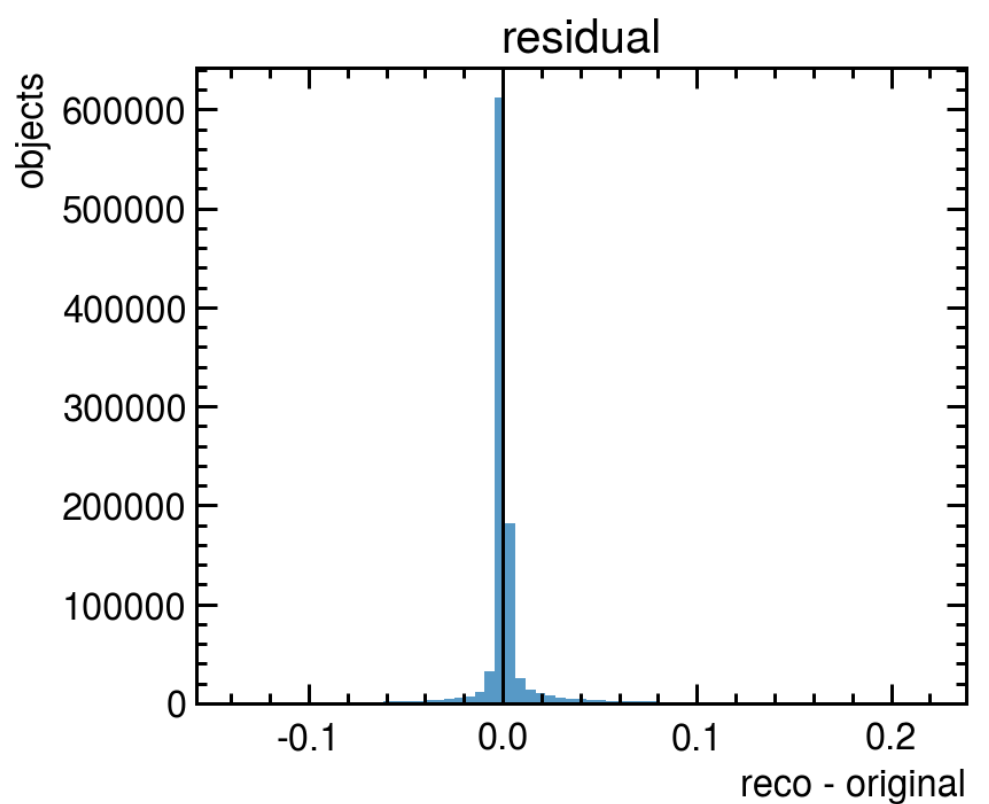
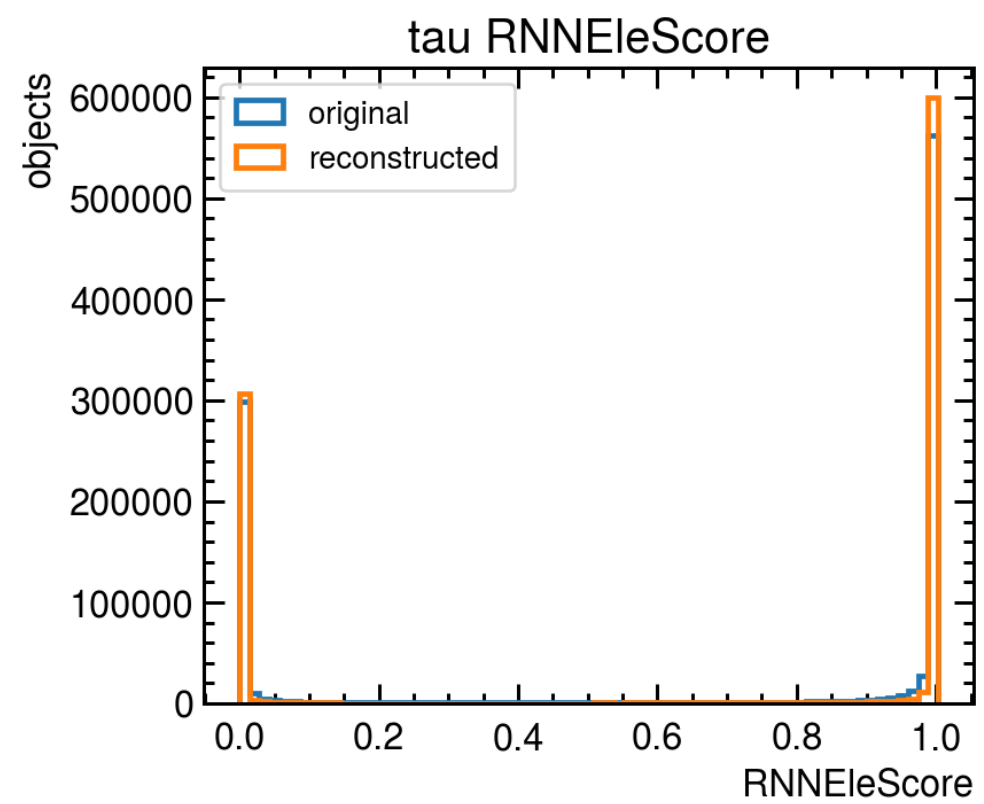
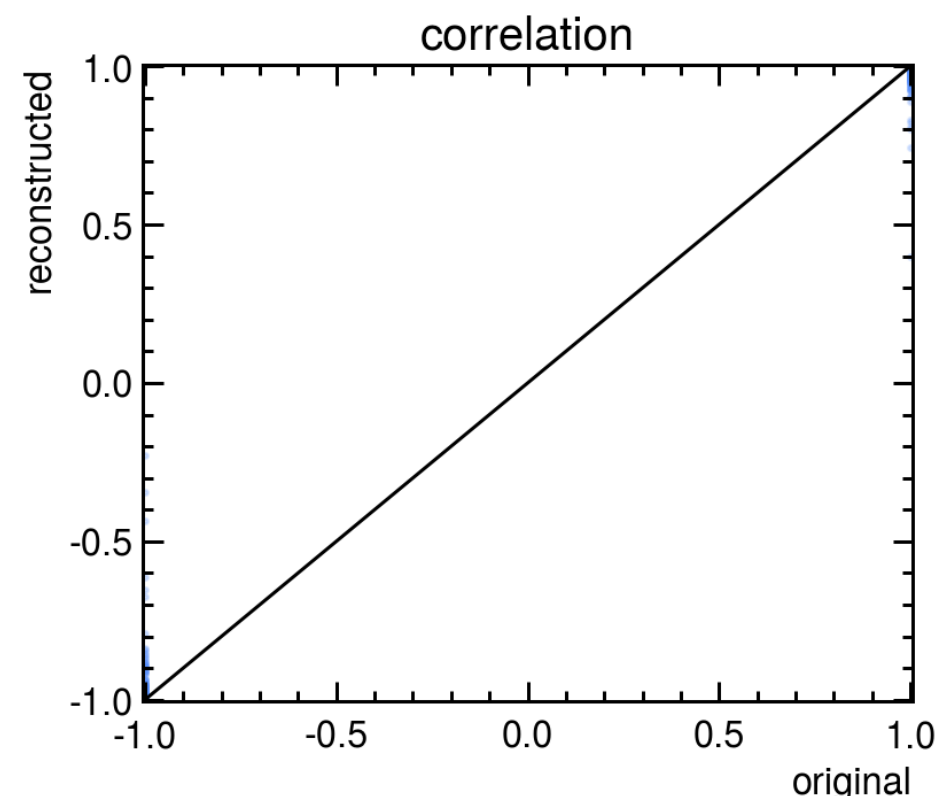
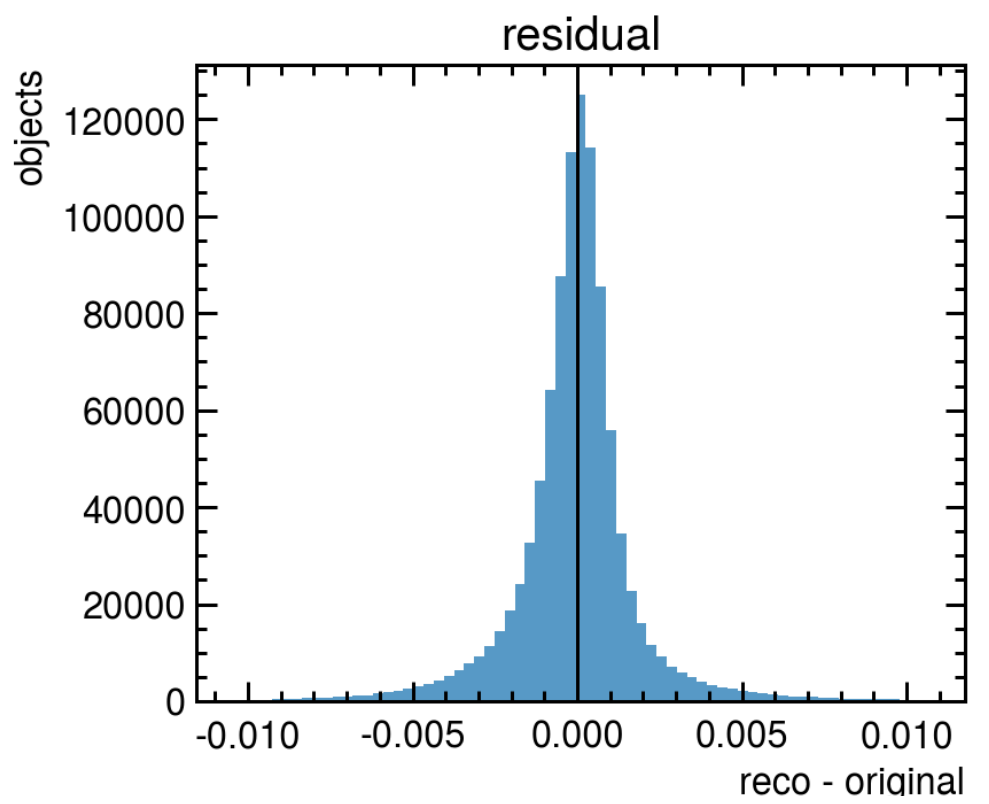
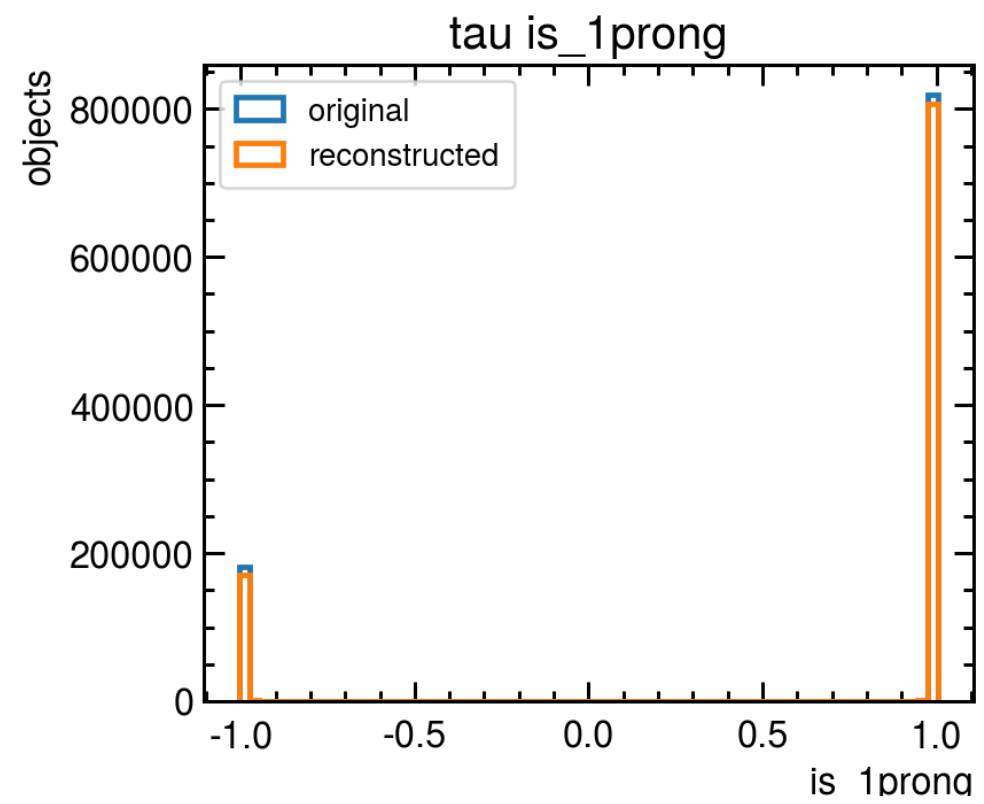
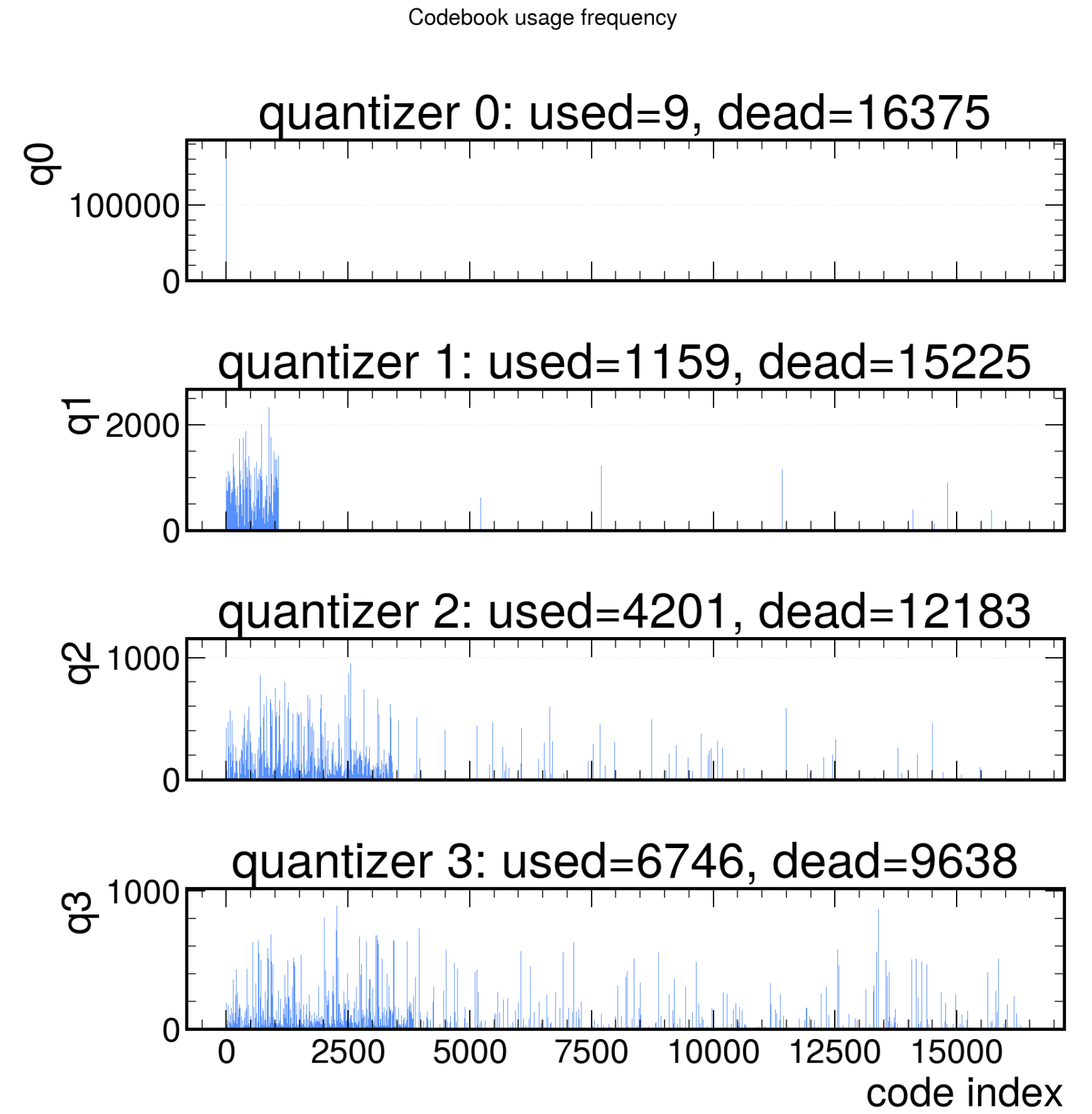
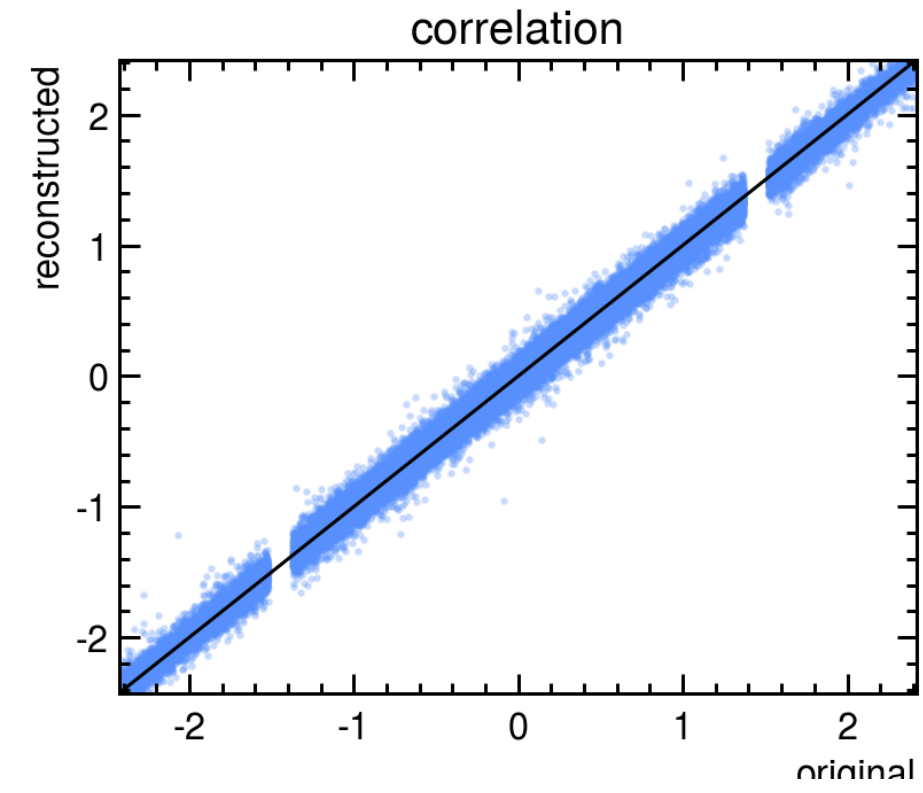
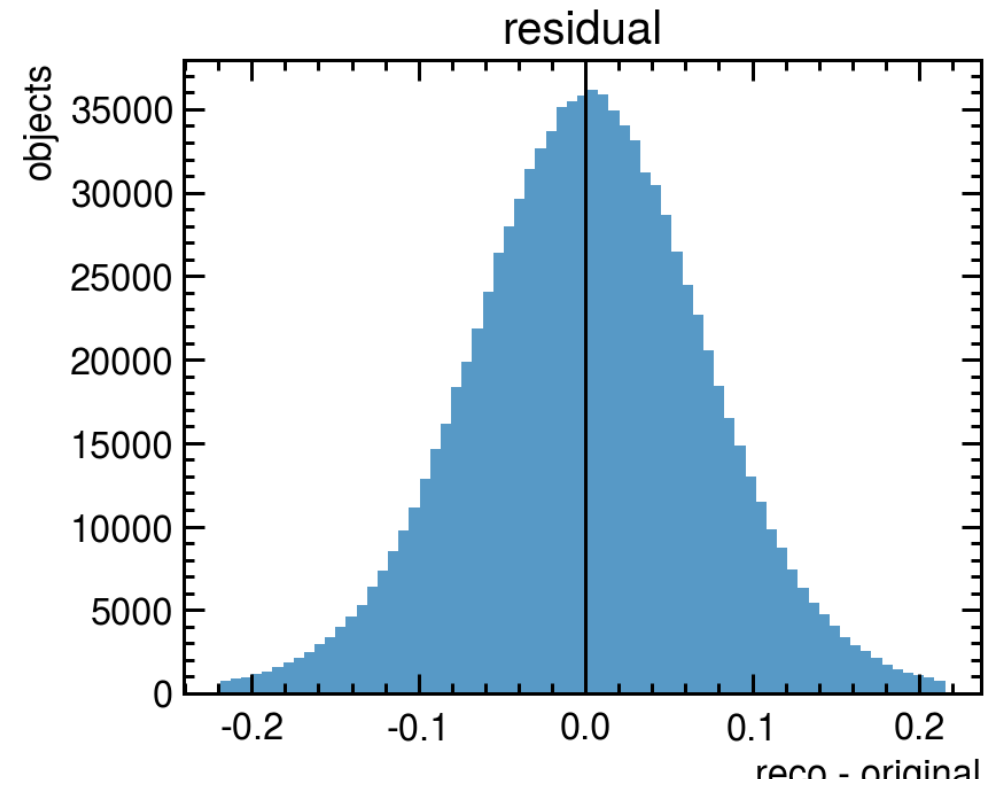
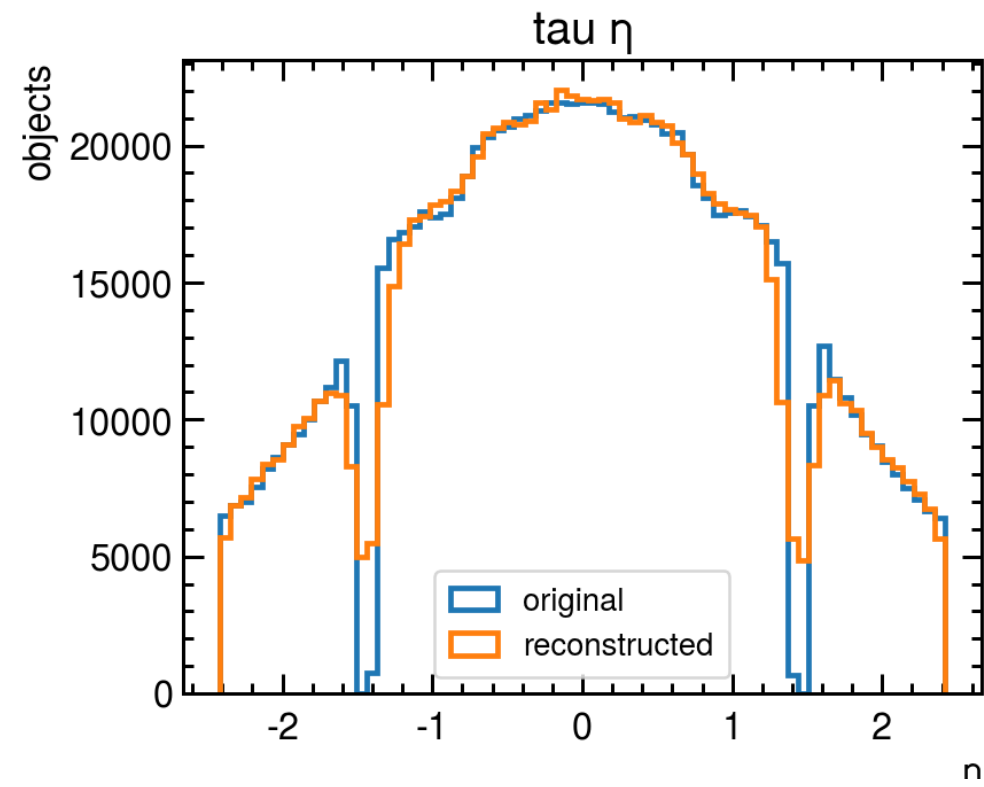
Reco - Jets



- * The best reco: **taus_cb16384_q4**
- * Lowest val/total_loss, val/recon_loss
- * Best reconstruction for pt, eta, phi
- * Clear improvement for RNNjetScore
- * An active fourth quantizer, with about 950 used q3 codes and increasing perplexity



Reco



* Summary

- ▶ 16384 x4 gives the best reconstruction for electrons, jets, and taus.
- ▶ **Two issues: q0 collapse, and many unused codebook entries.**
 - **Encoder outputs occupy a small number of coarse clusters.**
 - **q0 only needs a few codes to describe the broad structure.**
- ▶ The fourth quantizer is active and improves reconstruction.
- ▶ Larger codebooks improve some difficult features, but much of the capacity remains unused.
- ▶ Dead codes show inefficient use of capacity, but do not necessarily mean poor reconstruction.

* Possible Reasons

- ▶ The feature space may only contain a small number of coarse object patterns
- ▶ Continuous, binary, and categorical variables are all trained with the same reconstruction loss.
- ▶ Some features have very different scales or contain little useful variation.
- ▶ Codebook initialization, commitment loss, or encoder output scale may contribute to q0 collapse.
- ▶ Data and simulation differences may affect reconstruction.

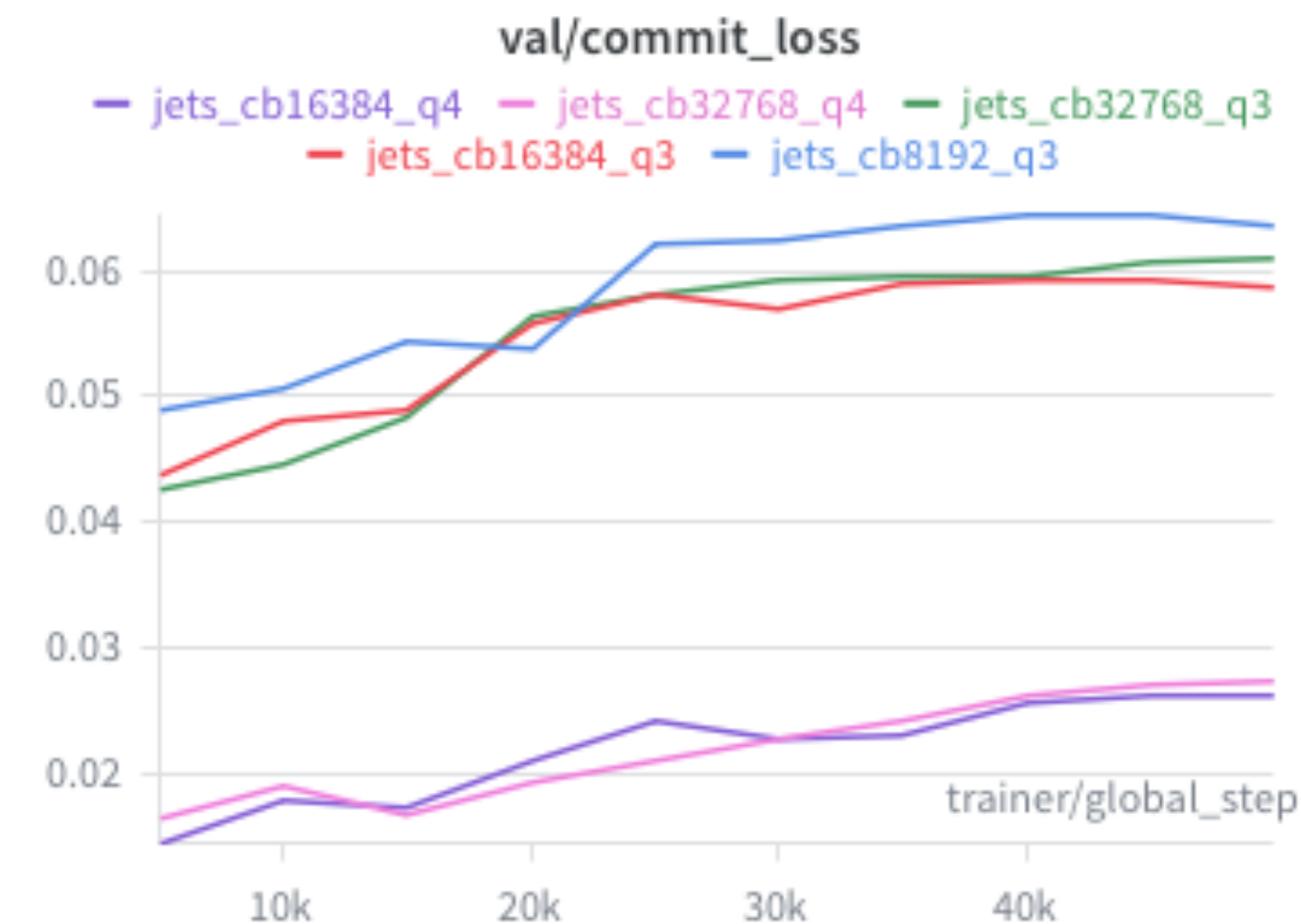
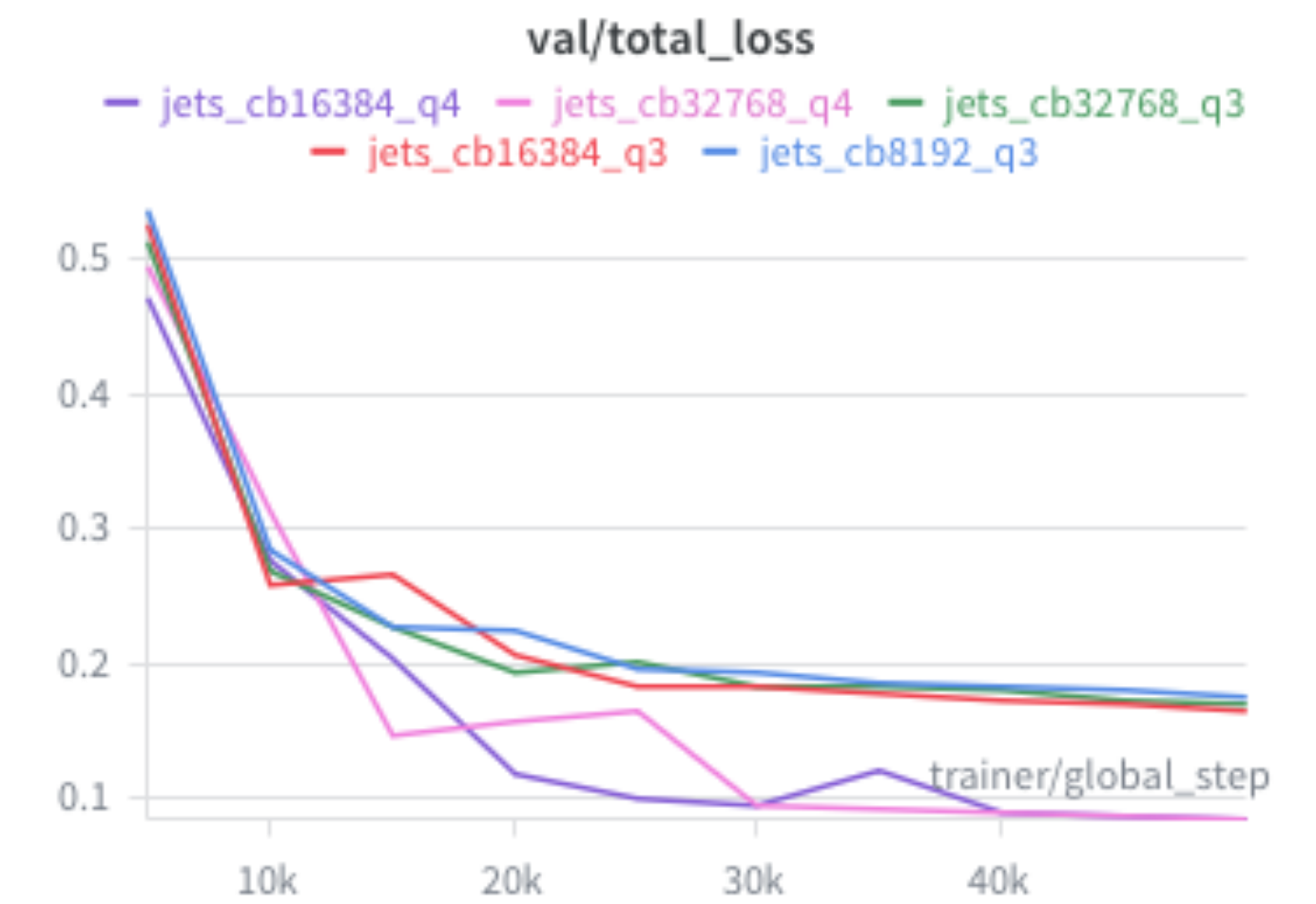
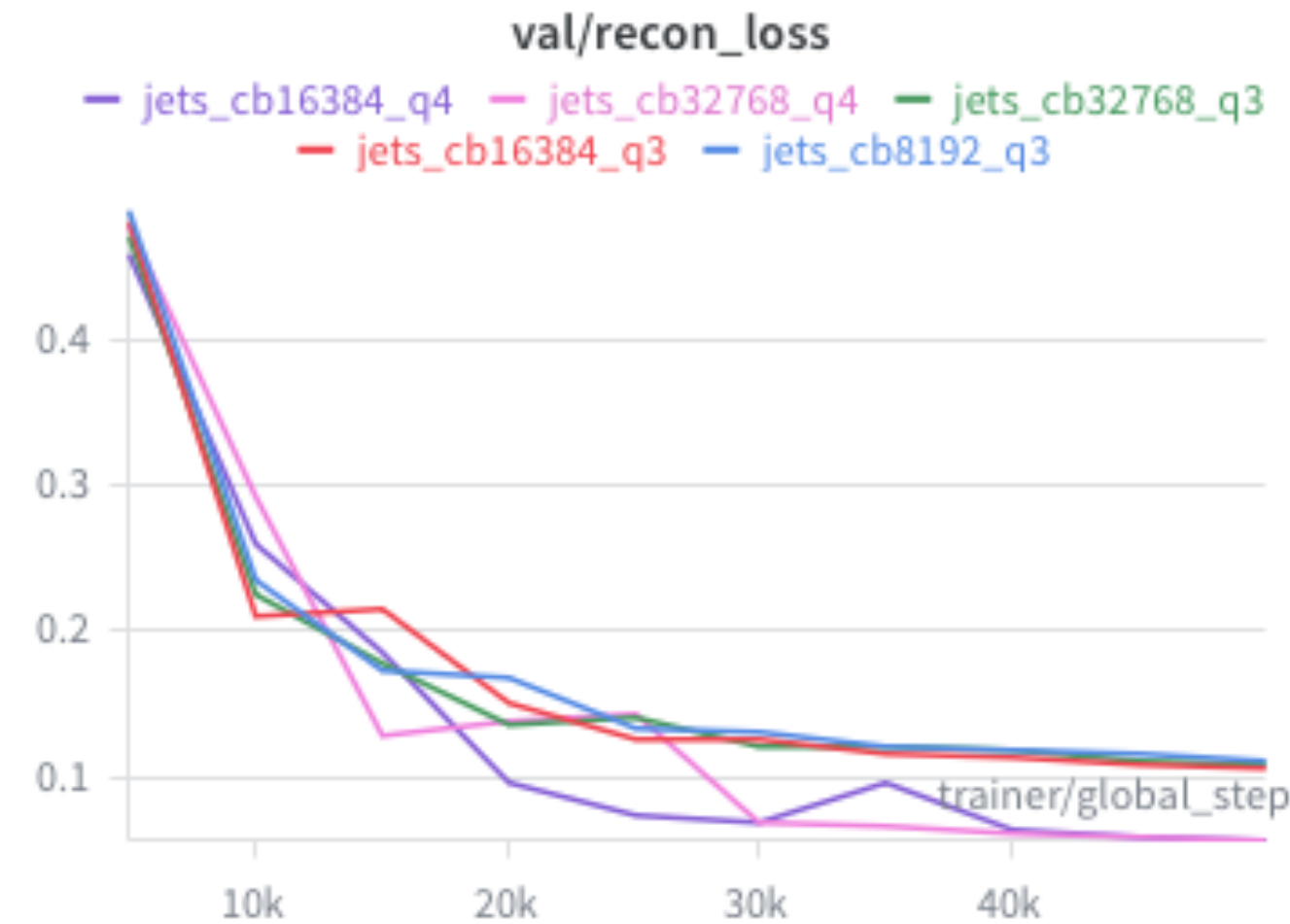
Next things to try

- * Try feature weighting or separate losses for continuous, tagger, binary, and categorical variables.
- * Scan smaller codebook_dim values, such as 8.
- * Add k-means codebook initialization instead of random (q0 check)
- * Reset dead codes during training (Codes used fewer than the threshold are replaced with vectors from the current batch)
- * Try cosine-similarity quantization
- * Log the encoder output before changing more parameters
- * Compare tokenizers using downstream transformer
- * Add derived physics checks such as invariant masses and angular separations after decoding



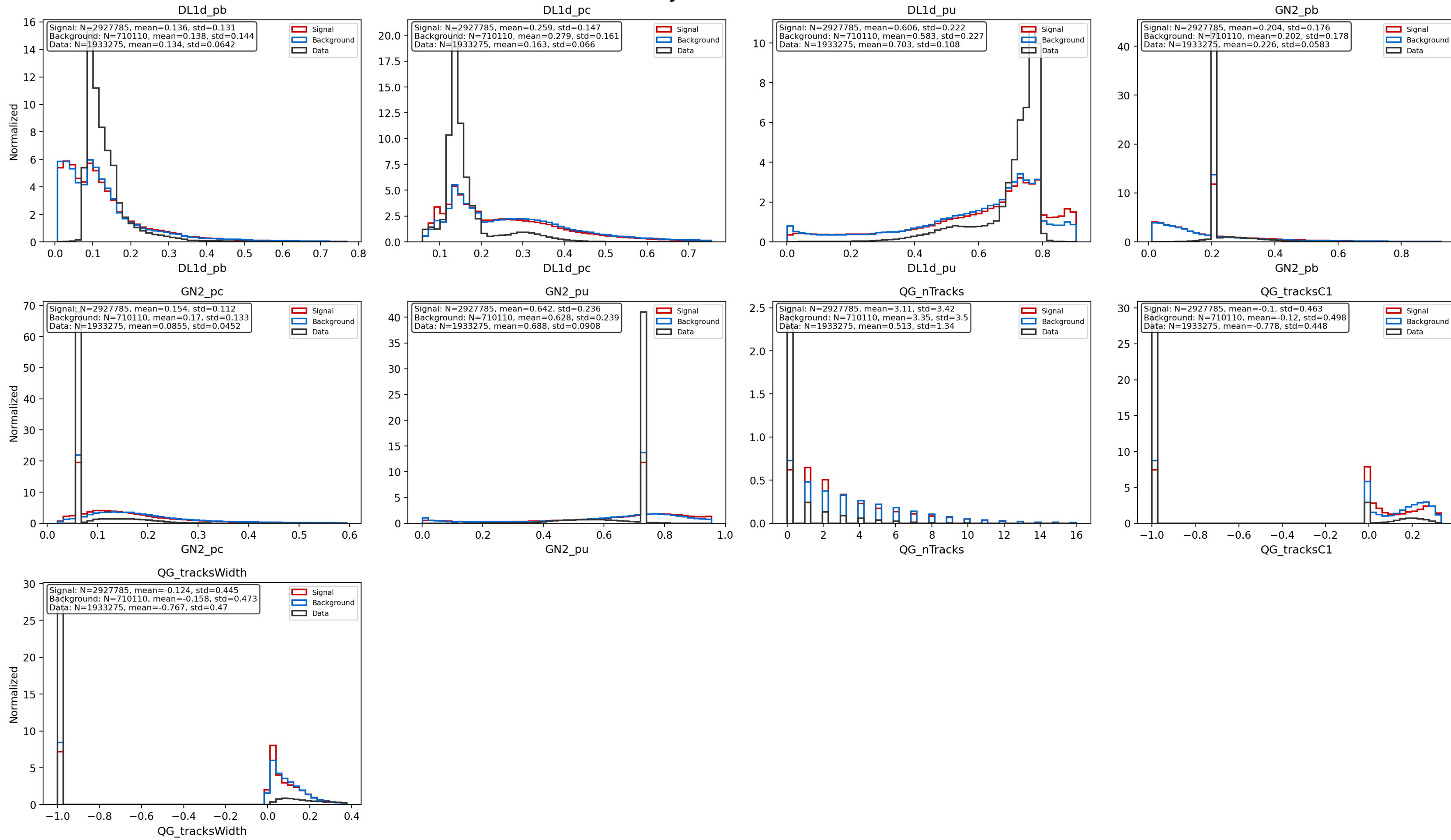
Loss

- * Commitment loss encourages encoder outputs to stay close to selected code vectors
- * $\text{total_loss} = \text{reconstruction_loss} + \text{commitment_loss}$
- * $\text{commitment_weight} = 1.0$
- * Commitment weight is contributing to optimisation but still smaller
- * Try different weights?
- * Try different initialisation of the codebooks?

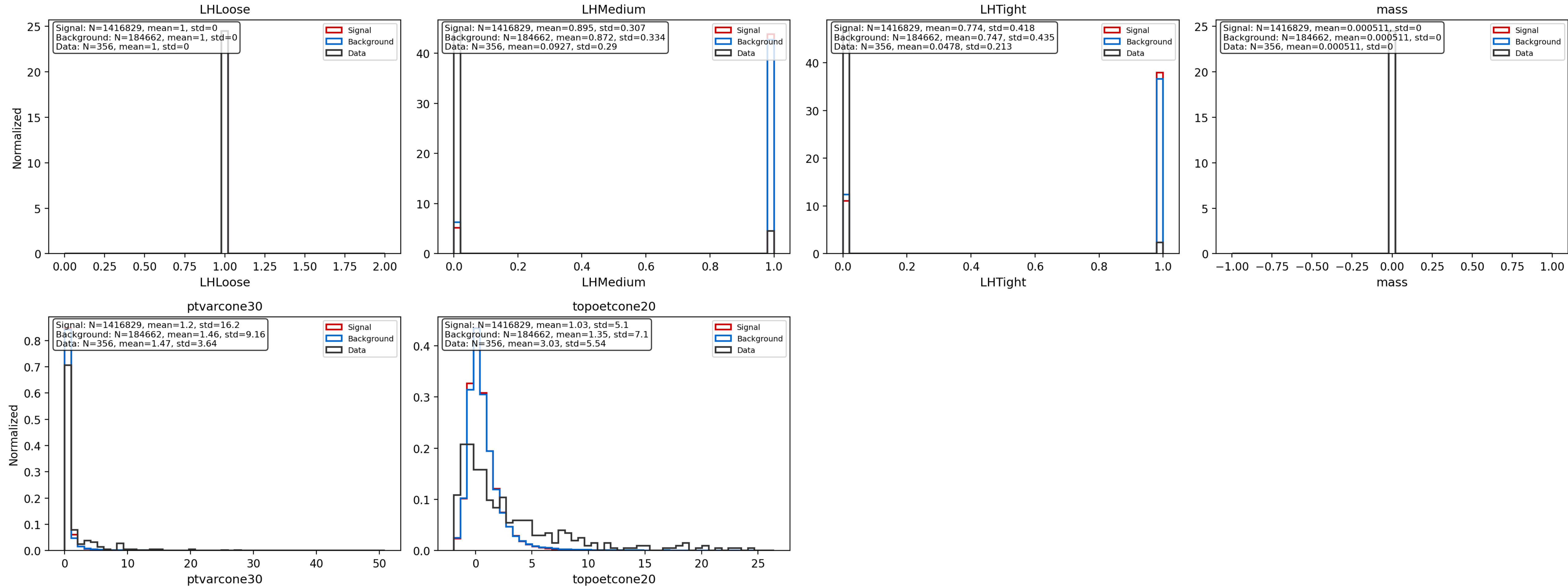


Data vs MC

atlas/jets/ — all features



atlas/electrons/ — all features



common/event/ — all features

