

Skeptical review: Predicting Halo Mass Function Proxies from Merger Tree Distributions using a Hybrid GNN and Gaussian Mixture Model

Summary

This manuscript proposes a hybrid ML pipeline to learn from ensembles of dark-matter halo merger trees and to predict, for each tree, a 20-bin histogram of (preprocessed) halo masses that the authors call an “HMF proxy” (Sec. 1, Sec. 2.5, Sec. 3.3). Each merger tree is represented as a graph with node features (mass, concentration, v_{\max} , scale factor) and encoded by a 3-layer GCN with global mean pooling into a 32D graph embedding (Sec. 2.1–2.3). A Gaussian Mixture Model is then fit to the embedding distribution and the per-tree posterior responsibilities are used as inputs to a small FFNN that predicts the per-tree mass histogram (Sec. 2.4–2.6). The paper also includes a pretext task in which embeddings are used to regress (Ω_m, σ_8) (Sec. 3.1), reports BIC-based selection of 5 GMM components with PCA/t-SNE visualizations (Sec. 3.1–3.2), and achieves a low reported MSE on the histogram target (Sec. 3.3). The overall direction—combining learned graph representations with mixture modeling to produce compact population descriptors—is promising and could be useful for emulation and conditional population modeling. However, several conceptual mismatches (especially “HMF” vs. within-tree histograms), under-specified data/simulation context, and unclear optimization details around the non-differentiable GMM currently limit scientific interpretability and reproducibility. In addition, because mass is both a primary node feature and the prediction target, stronger baselines/ablations are needed to establish that the method learns more than a compressive “mass-in/mass-out” mapping and to justify the added complexity of the GMM stage.

Strengths

- Addresses a relevant cosmology problem: learning informative representations of merger trees to support downstream prediction tasks related to halo populations (Sec. 1).
- Clear modular design (GNN \rightarrow embedding distribution modeling via GMM \rightarrow FFNN prediction), which is in principle interpretable and extensible (Sec. 2.1–2.6).
- Includes a physically motivated auxiliary/pretext evaluation (cosmological parameter regression) suggesting embeddings capture nontrivial information (Sec. 3.1).
- Uses standard tools (BIC, PCA/t-SNE) to probe embedding structure and clustering behavior (Sec. 3.1–3.2).
- Provides qualitative predicted-vs-true histogram comparisons and a quantitative test metric for the main task (Sec. 3.3; Fig. 4/5 as referenced).
- Figures are generally readable and the paper is largely well organized at a high level, making the intended workflow easy to follow.

Major issues

1. **Conceptual mismatch and insufficient definition of the target: the paper frames the prediction target as an “HMF proxy”, but it appears to be a within-tree halo/progenitor mass histogram (often described as normalized) rather than a standard halo mass function (number density per comoving volume at a given redshift) (Sec. 1, Sec. 2.5, Sec. 3.3, Sec. 3.5, Sec. 4). Relatedly, it is unclear at what times/redshifts the masses entering the histogram are taken (all nodes across all snapshots? only a specific scale factor? only progenitors above a threshold), which fundamentally changes the physical meaning of the distribution.**

Recommendation: In Sec. 1 and Sec. 2.5, (i) formally define the conventional HMF (e.g., $dn/d\log M$ per comoving volume at fixed redshift) and (ii) precisely define the paper’s target with an explicit formula. Specify exactly which nodes contribute to the histogram (all halos across the full tree vs. only at a chosen snapshot; inclusion/exclusion of subhalos; any mass thresholds) and the redshift/scale-factor convention. To avoid overclaiming, rename the target throughout to something like “per-tree progenitor mass distribution / within-tree mass spectrum” unless you provide a quantitative mapping to an HMF. If the goal is ultimately HMF emulation, add a concrete demonstration in Sec. 3.3 or Sec. 3.5 that aggregates predictions across many trees (with correct weighting/normalization, and specifying the effective volume and selection function) and compares to an HMF measured directly from the simulation.

2. **Potential leakage / ill-posedness of the main task: node features include (log) mass, and the target is a histogram of masses from the same tree. As written, the pipeline may succeed primarily by compressing and reconstructing the mass distribution (a near-autoencoding problem) rather than learning meaningful formation-history/topology-to-population relations (Sec. 2.2–2.5, Sec. 3.3). Without strong baselines/ablations, it is difficult to assess what is learned beyond “mass-in, mass-out.”**

Recommendation: Add explicit baselines and ablations in Sec. 3.3: (i) a non-graph baseline that predicts the histogram from simple mass-only summaries (moments/quantiles, counts above thresholds, or even the raw list of masses pooled via a DeepSets-style model); (ii) an ablation removing the mass feature from node inputs (using only concentration, v_{\max} , scale factor, and structure) to test whether non-mass information contributes; and (iii) an ablation using only mass (and optionally scale factor) to quantify how much topology/other features help. Interpret the results accordingly in Sec. 3.5/4: if mass-only is near-optimal, the claim should shift from “formation-history inference” to “learned compression of mass spectra.”

3. **Inconsistent and under-specified mass preprocessing and binning, including a contradiction between the single-log transform in Sec. 2.1 (Eq. (1)) and the double-log statement in Sec. 3.4 (“ $\log_{10}(\log_{10}(\text{Mass}))$ ”), as well as ambi-**

guity about log base, offsets, and whether standardized values are being conflated with “denormalized” ones (Sec. 2.1, Sec. 2.5, Sec. 3.4; Fig. 5 as referenced). This undermines interpretability of the reported mass ranges and the predicted histograms.

Recommendation: Unify the preprocessing description across Sec. 2.1, Sec. 2.5, and Sec. 3.4 by giving one explicit, end-to-end definition of the mass variable used for (a) node features and (b) histogram binning/targets. State: (i) whether you use \ln or \log_{10} ; (ii) whether there is one log or two; (iii) any dimensionless scaling (e.g., $\log_{10}(M/M_0 + \epsilon)$ with M_0 and ϵ specified in the same units as M); and (iv) where standardization is applied and how it is inverted for plotting. Provide a table of the 20 bin edges both in transformed space and in physical mass units, and ensure figure axis labels/captions match the actual transform. Correct Sec. 3.4 if the double-log is a typo/artifact; if a double-log is genuinely used, justify it and discuss its implications for physical interpretation.

4. **Training protocol and the role of the GMM in the claimed “joint” GNN–GMM–FFNN optimization are unclear and may be technically inconsistent (Sec. 2.4–2.6). The forward pass uses GMM posterior responsibilities, but standard sklearn `GaussianMixture` fitting is non-differentiable in PyTorch, making it unclear whether gradients can propagate through the GMM to the GNN, whether the GMM is frozen, and whether embeddings drift away from the fitted mixture during training.**

Recommendation: In Sec. 2.4 and Sec. 2.6, provide an explicit step-by-step training schedule (preferably pseudocode/Algorithm box): (1) whether/how the GNN is pre-trained (e.g., on the cosmology task in Sec. 3.1) and whether it is frozen; (2) when the GMM is fit (once on training embeddings vs. refit periodically/alternating optimization); (3) whether GMM parameters are held fixed; and (4) whether responsibilities are treated as constants for backprop. If the GMM is fit once and frozen, rephrase “joint training” accordingly and discuss embedding drift; if you refit GMM during training, describe the alternation frequency and convergence criteria. If you implemented a differentiable mixture model in torch, describe it and cite/justify it.

5. **Insufficient description of simulations, cosmology coverage, and dataset construction; possible train/test leakage due to non-independence of trees (Sec. 2.1–2.2, Sec. 3.1, Sec. 3.4). The manuscript does not specify the simulation code, box size/volume, mass resolution, halo finder and mass definition (e.g., M_{200c}), tree-builder, redshift range/snapshots, or how the 1000 trees are selected (by root mass? random halos? hosts/subhalos?). It is also unclear whether trees span multiple cosmologies and how (Ω_m, σ_8) labels are assigned per tree, and whether the split avoids leakage across cosmologies or correlated halos in the same volume.**

Recommendation: Add a dedicated dataset/simulation subsection (expand Sec. 2.1) containing: simulation name/code, box size (volume), particle mass and force resolution, halo finder, halo mass definition, tree-building algorithm, snapshot/redshift coverage, and units. Describe tree selection criteria (root halo mass cuts, redshift of root, host/subhalo handling) and provide summary statistics of graph sizes (#nodes/#edges) and root masses. Clarify whether the dataset includes multiple cosmologies; if so, state how many and the parameter ranges, and ensure splits are performed by cosmology (or by simulation realization) to test generalization rather than memorization. If all trees come from one volume, discuss correlation risks and consider split strategies that reduce shared-mode leakage (e.g., by spatial region, by halo mass bins, or by simulation realizations if available).

6. **Target construction, output constraints, and loss/metrics are not mathematically aligned: the target is described as a normalized histogram/probabilities (Sec. 2.5, Sec. 3.3), but the FFNN uses independent sigmoids per bin, which does not enforce sum-to-one, and MSE alone is hard to interpret for histogram comparisons (Sec. 2.5–2.7).**

Recommendation: First, state explicitly in Sec. 2.5 whether the target is (a) counts per bin or (b) a normalized probability mass function. If (b), enforce normalization (softmax over 20 bins or explicit L1 normalization) and use a distribution-aware loss (cross-entropy/KL/JS; optionally add EMD/Wasserstein as an evaluation metric). If (a), remove probability language, justify sigmoid bounds (or use nonnegative outputs such as softplus), and evaluate with metrics appropriate for counts (possibly Poisson/negative-binomial likelihood). In Sec. 3.3, supplement global MSE with per-bin error plots (especially high-mass tail), and at least one distributional metric (KL/JS or EMD) to demonstrate shape fidelity beyond average squared error.

7. **Insufficient evidence that the specific hybrid design (GNN \rightarrow GMM responsibilities \rightarrow FFNN) is necessary: the evaluation reports a single headline test MSE on 100 trees with limited diagnostics, no uncertainty across random seeds/splits, and no ablation demonstrating the value of the GMM step versus directly using embeddings (Sec. 3.3, Sec. 3.5).**

Recommendation: In Sec. 3.3 (or a new ablation subsection), add: (i) direct GNN-embedding \rightarrow FFNN prediction (no GMM) and compare performance; (ii) alternative clustering/bottlenecks (e.g., k-means responsibilities, vector quantization) to test whether the benefit is “mixture modeling” specifically; (iii) multiple runs with different seeds and/or multiple data splits, reporting mean \pm std of key metrics; and (iv) stratified performance by root mass, tree size, and (if kept) GMM component. This will quantify what the GMM contributes and reveal systematic failure modes.

8. **GMM clustering analysis is currently qualitative and not tied to physical/structural differences among merger trees (Sec. 3.2). BIC selection is presented on a very sparse grid in Fig. 2 and cluster meaning is inferred**

primarily from low-dimensional projections, which can be misleading without quantitative characterization.

Recommendation: Strengthen Sec. 3.2 by: (i) scanning a contiguous range of component counts (e.g., $K = 1-20$) with multiple initializations and reporting variability; (ii) specifying which split (train only) was used to fit/score the GMM; (iii) reporting mixture weights and quantitative cluster separation (e.g., silhouette on embeddings; or within-/between-component covariance diagnostics); and (iv) summarizing each component with physically interpretable tree statistics (root mass, node count, formation time proxies like last major merger scale factor, concentration/ v_{\max} distributions, etc.). If components do not map cleanly to physical categories, state that explicitly as a limitation.

- 9. The cosmology pretext task is under-specified and its scientific meaning is unclear without target ranges/normalization, per-parameter errors, and proper baselines; it is also unclear whether the evaluation tests generalization across cosmologies or merely across trees within the same cosmology (Sec. 3.1, Sec. 4).**

Recommendation: In Sec. 3.1, specify: the number of cosmologies, parameter ranges for (Ω_m, σ_8) , whether targets are standardized, the regression head architecture, and the exact metric definition. Report per-parameter RMSE/MAE in physical units (not only combined MSE), include predicted-vs-true scatter plots, and compare against simple baselines (predict mean; linear regression on hand-crafted tree summaries). If multiple cosmologies exist, evaluate with a split by cosmology (hold out cosmologies) to support the claim that embeddings capture cosmology information rather than overfitting to particular simulations/realizations.

Minor issues

1. Merger trees are inherently directed and time-ordered, but the current GCN description does not clarify whether edge direction and temporal information are used or discarded; no edge features are described (Sec. 2.1–2.3). Ignoring direction/time may limit physical expressiveness and interpretability.

Recommendation: Clarify in Sec. 2.1–2.3 whether the graph is treated as directed or undirected in message passing and how temporal ordering is represented (if at all). If direction/time is ignored, briefly justify; otherwise consider adding edge attributes (time gaps, merger mass ratios) or using directed/temporal message passing, and include an ablation showing whether it improves results.

2. Exploratory data analysis (EDA) is referenced to justify choices like the histogram range and binning but is not shown quantitatively (Sec. 2.2, Sec. 3.4).

Recommendation: Add a compact EDA summary: distribution plots and key statistics (min/median/max, percentiles) of the mass variable used for binning, plus typical node counts per tree. Use these to justify the chosen mass range and the choice of 20 bins (Sec. 2.5).

3. t-SNE plots risk overinterpretation and are missing reproducibility details; they also appear to be shown only for training embeddings and without coloring by meaningful physical attributes (Sec. 3.1; Fig. 1).

Recommendation: Provide t-SNE hyperparameters (perplexity, learning rate, iterations, seed) in the caption or methods, show stability across a couple of settings/seeds, and avoid strong claims about global structure. Consider coloring by root mass, tree size, cosmology label, or GMM component responsibilities, and include validation/test embeddings to check out-of-sample consistency.

4. BIC/AIC model selection for the GMM is presented only for $K \in 5, 10, 15$ and without variability across initializations; the selected K is not visually emphasized (Sec. 3.2; Fig. 2).

Recommendation: Expand the sweep to a denser range (e.g., 1–20) and run multiple initializations to show variability. Indicate clearly which K is selected and by what margin (ΔBIC), and state whether scoring is on training only.

5. Related work and novelty positioning are thin: the introduction does not sufficiently situate the method relative to existing HMF emulators, conditional mass function/assembly bias modeling, and prior ML/GNN work on merger trees (Sec. 1).

Recommendation: Add a brief related-work subsection near the end of Sec. 1 summarizing prior HMF emulation and ML on merger trees/GNNs in cosmology, and then state clearly what is novel here (e.g., using mixture-model responsibilities as distribution-level features for downstream prediction).

6. Presentation: section formatting is inconsistent (e.g., mixed heading markers with “#”), some keywords are off-topic, and several figures would benefit from higher-resolution export and more explicit caption details (Sec. 2.7, Sec. 3; Title/Keywords; Figs. 1–4).

Recommendation: Standardize section headings/numbering, revise keywords to match the actual content (halo mass function, dark matter halos, merger trees, GNNs, emulation), and improve figure export (vector/high DPI), font sizes, and captions (sample sizes, normalization, bin counts).

Very minor issues

1. Notation/typography inconsistencies: v_{\max} vs v_{max} , “log” vs “log10”, and Eq. (3) uses scalar notation although standardization is feature-wise (Sec. 2.1, Sec. 3.4).

Recommendation: Standardize notation throughout (e.g., v_{\max}) and specify the log base everywhere. Add a short note after Eq. (3) clarifying that μ and σ are computed per feature over the training set.

2. Some repeated phrasing in Results/Conclusion and mildly informal wording reduce concision (Sec. 3.3–3.5, Sec. 4).

Recommendation: Edit for concision by removing near-duplicate sentences and tightening the Conclusion to emphasize (i) the main quantitative findings, (ii) the clarified target definition, and (iii) concrete limitations and next steps.

Mathematical consistency audit

This section audits **symbolic/analytic** mathematical consistency (algebra, derivations, dimensional/unit checks, definition consistency).

Maths relevance: light

The paper contains only a small number of explicit equations (log transforms and z-score standardization), plus several mathematically loaded statements about probability-normalized histograms and the use of GMM posterior probabilities in a joint-training pipeline. The main internal-consistency concerns are (i) conflicting definitions of the transformed mass feature (single log vs double log₁₀), (ii) interpreting sigmoid outputs as a probability histogram without enforcing normalization, and (iii) an underspecified optimization/gradient-flow relationship between the GNN, a separately fit GMM, and the FFNN.

Checked items

1. ✘ **Mass log-transform definition** (Eq. (1), Sec. 2.1, p.2)
 - **Claim:** Mass feature is transformed as $\text{masstransformed} = \log(\text{mass} + 10^{-6})$.
 - **Checks:** symbol/notation consistency, dimensional consistency, cross-section consistency
 - **Verdict:** FAIL; confidence: high; impact: critical
 - **Assumptions/inputs:** “log” denotes some logarithm (base not specified), mass is nonnegative; 10^{-6} prevents $\log(0)$.
 - **Notes:** This definition conflicts with Sec. 3.4 (p.6), which claims the effective mass feature used is $\log_{10}(\log_{10}(\text{Mass}))$. No derivation/step is provided linking these, so the mass feature used for binning/plots is ambiguous and internally inconsistent.
2. ✔ **Concentration log-transform definition** (Eq. (2), Sec. 2.1, p.2)
 - **Claim:** Concentration feature is transformed as $\text{concentrationtransformed} = \log(\text{concentration} + 10^{-6})$.
 - **Checks:** algebraic correctness, domain/sanity
 - **Verdict:** PASS; confidence: medium; impact: minor

- **Assumptions/inputs:** Concentration is dimensionless and nonnegative., 10^{-6} is a small offset in the same units (dimensionless here).
 - **Notes:** Algebraically fine. Base of log is unspecified (minor clarity issue).
3. ✓ **Z-score standardization formula** (Eq. (3), Sec. 2.1, p.2)
- **Claim:** Standardization is $x_{\text{scaled}} = (x - \mu)/\sigma$ using training-set mean and standard deviation.
 - **Checks:** algebraic correctness, definition consistency
 - **Verdict:** PASS; confidence: high; impact: minor
 - **Assumptions/inputs:** $\sigma \neq 0$ for each standardized feature.
 - **Notes:** Standard and correct. Would benefit from stating explicitly that μ and σ are computed per feature.
4. ✓ **Scale-factor normalization range claim** (Sec. 2.1, bullet 2, p.2)
- **Claim:** Scale factor is normalized to $[0, 1]$ by dividing by the maximum scale factor observed in the dataset.
 - **Checks:** sanity/limiting case
 - **Verdict:** PASS; confidence: medium; impact: minor
 - **Assumptions/inputs:** Scale factor values are nonnegative.
 - **Notes:** Dividing by the maximum yields values in $[0, 1]$ provided the maximum is positive; typical scale factor values satisfy this.
5. ✗ **HMF proxy as normalized histogram vs sigmoid outputs** (Sec. 2.5, p.3; Sec. 3.3, p.5)
- **Claim:** The HMF proxy is a normalized histogram with 20 bins (probabilities), and the FFNN output uses a sigmoid so predicted values lie in $[0, 1]$ and represent probabilities.
 - **Checks:** normalization/constraints, definition consistency
 - **Verdict:** FAIL; confidence: high; impact: critical
 - **Assumptions/inputs:** “Normalized histogram” means bin values sum to 1 for each tree.
 - **Notes:** A sigmoid applied independently to 20 outputs only constrains each bin to $(0, 1)$ and does not enforce sum-to-1. This is inconsistent with a probability distribution over bins unless an additional normalization step (not stated) is applied.
6. ✓ **GMM posterior probabilities as FFNN input** (Sec. 2.5, p.3)
- **Claim:** FFNN input is the vector of GMM posterior probabilities for each tree (length = number of components).
 - **Checks:** constraint/sanity (sum-to-one)
 - **Verdict:** PASS; confidence: medium; impact: minor

- **Assumptions/inputs:** Posterior responsibilities are computed correctly by the GMM.
 - **Notes:** A posterior-probability vector is well-defined and typically sums to 1 across components, making it a coherent feature representation.
7. \triangle **Joint training with intervening GMM step** (Sec. 2.6, steps 1–6, p.3–4)
- **Claim:** The GNN and FFNN are trained jointly via backprop on MSE loss, with the GMM used to compute posterior probabilities from embeddings inside the forward pass.
 - **Checks:** derivation/optimization logic consistency, missing steps
 - **Verdict:** UNCERTAIN; confidence: medium; impact: critical
 - **Assumptions/inputs:** GMM parameters are either fixed during training or updated in a specified alternating scheme., Posterior probabilities are treated as differentiable functions of embeddings if gradients are claimed w.r.t. GNN parameters.
 - **Notes:** The paper does not specify whether the GMM is fixed, refit, or jointly optimized while the GNN changes embeddings. Without that, it is not possible to verify that the stated joint backpropagation procedure is mathematically coherent as described.
8. \checkmark **Mapping from $\log_{10}(\log_{10}(\text{Mass}))$ range to mass range** (Sec. 3.4, p.6)
- **Claim:** If $y = \log_{10}(\log_{10}(\text{Mass}))$ ranges from 0.9829 to 1.1619, then **Mass** ranges from $10^{9.61}$ to $10^{14.52}$ (in stated units).
 - **Checks:** algebraic inversion
 - **Verdict:** PASS; confidence: high; impact: moderate
 - **Assumptions/inputs:** Definition $y = \log_{10}(\log_{10}(M))$ uses base-10 logs., M is expressed in consistent positive units.
 - **Notes:** Given $y = \log_{10}(\log_{10}(M))$, then $\log_{10}(M) = 10^y$ and $M = 10^{10^y}$. The reported mapping corresponds to $10^y \approx 9.61$ and 14.52 , hence $M \approx 10^{9.61}$ to $10^{14.52}$, which is consistent. However, this conflicts with Eq. (1) elsewhere.
9. \triangle **Dimensional validity of $\log(\text{mass} + 10^{-6})$** (Eq. (1), Sec. 2.1, p.2)
- **Claim:** Taking log of mass after adding 10^{-6} is a valid preprocessing step.
 - **Checks:** dimensional/units consistency
 - **Verdict:** UNCERTAIN; confidence: low; impact: moderate
 - **Assumptions/inputs:** mass carries physical units.
 - **Notes:** Logarithms require dimensionless arguments. The paper does not state whether mass was nondimensionalized (e.g., divided by a reference mass) before applying the log, nor the units context for the additive 10^{-6} . This is a conceptual dimensional ambiguity.

Limitations

- Only the provided 7-page PDF text/images were used; there are no extended derivations or appendices to audit.
- The audit does not verify any reported numerical values, simulation results, plots, or implementation feasibility beyond analytic/logical consistency.
- Because the paper omits explicit formulas for the GMM posterior (responsibilities) and the exact histogram normalization procedure, some checks (especially around training/joint optimization and target normalization) are necessarily based on stated claims rather than fully specified equations.

Numerical results audit

This section audits **numerical/empirical** consistency: reported metrics, experimental design, baseline comparisons, statistical evidence, leakage risks, and reproducibility.

14/14 stated-value arithmetic/logical consistency checks passed. Key internal consistencies verified include dataset split fractions/counts, architecture dimensionality and output/bin alignment, derived MSE-per-bin, and mass-range back-calculation from double-log values.

Checked items

1. ✓ **C1_dataset_split_counts** (p.2, Sec. 2.1 Data Acquisition and Preprocessing)
 - **Claim:** “The dataset consisted of 1000 merger trees... split into training (80%), validation (10%), and testing (10%) sets.”
 - **Checks:** parts_vs_total / percent_to_count
 - **Verdict:** PASS
 - **Notes:** Fractions sum to 1 and implied integer counts match 800/100/100.
2. ✓ **C2_train_set_size_figure1_caption_vs_split** (p.4, Fig. 1 caption + p.2 Sec. 2.1 split description)
 - **Claim:** Fig. 1 caption: “32-dimensional GNN embeddings of the 800 training merger trees.” Split description implies 80% of 1000 = 800.
 - **Checks:** cross_reference_consistency (count)
 - **Verdict:** PASS
 - **Notes:** Caption count matches implied training size.
3. ✓ **C3_eps_constant_matches_equations** (p.2, Eqs. (1)-(2))
 - **Claim:** “A small constant of 10^{-6} was added...: $\log(\text{mass} + 10^{-6})...$ $\log(\text{concentration} + 10^{-6})$.”
 - **Checks:** repeated_constant_consistency
 - **Verdict:** PASS
 - **Notes:** Same constant used.

4. ✓ **C4_gcn_layer_dims_consistency** (p.3, Sec. 2.3 Graph Neural Network Architecture)
 - **Claim:** GCN layers: input→64, 64→32, 32→32; global mean pooling yields embedding size **32**.
 - **Checks:** dimension_chain_consistency
 - **Verdict:** PASS
 - **Notes:** Dimensions consistent.

5. ✓ **C5_hmf_bins_output_dim_match** (p.3, Sec. 2.5 Halo Mass Function Prediction)
 - **Claim:** “20 bins were used... Output Layer... 20 neurons (corresponding to the 20 bins in the HMF histogram).”
 - **Checks:** parameter_match (bins vs output units)
 - **Verdict:** PASS
 - **Notes:** Output units match bins.

6. ✓ **C6_gmm_component_options_count** (p.3, Sec. 2.4 Gaussian Mixture Model + p.5 Fig. 2 caption)
 - **Claim:** “Experiments... (5, 10, and 15) components...” and Fig. 2: “GMMs with 5, 10, and 15 components.”
 - **Checks:** cross_reference_consistency (list of tested values)
 - **Verdict:** PASS
 - **Notes:** Parsed set matches 5, 10, 15.

7. ✓ **C7_early_stopping_patience_vs_epochs** (p.4, Sec. 2.6 Training Procedure)
 - **Claim:** “Training... fixed number of epochs (100)... Early stopping... stopped if the validation loss did not improve for 10 consecutive epochs.”
 - **Checks:** simple_inequality_sanity
 - **Verdict:** PASS
 - **Notes:** Logical inequality holds.

8. ✓ **C8_optimizer_hyperparams_positive** (p.3-4, Sec. 2.6 Training Procedure)
 - **Claim:** “Adam optimizer... learning rate of 0.001 and a weight decay of 0.0001.”
 - **Checks:** range_check (non-negative hyperparameters)
 - **Verdict:** PASS
 - **Notes:** Hyperparameters satisfy non-negativity/positivity.

9. ✓ **C9_mse_per_bin_from_total_mse** (p.5, Sec. 3.3 HMF Proxy Prediction)
 - **Claim:** “FFNN achieved a mean MSE of 0.000522... Given... 20 bins, the average MSE per bin is approximately 2.61×10^{-5} .”

- **Checks:** derived_value (division)
 - **Verdict:** PASS
 - **Notes:** Claim matches division within tolerance.
10. ✓ **C10_mse_std_less_than_mean_check** (p.5, Sec. 3.3 HMF Proxy Prediction)
- **Claim:** “mean MSE of 0.000522... standard deviation of 0.000386.”
 - **Checks:** basic_sanity (dispersion vs mean)
 - **Verdict:** PASS
 - **Notes:** Std is nonnegative and less than mean.
11. ✓ **C11_pca_variance_explained_bounds** (p.5, Sec. 3.2 + Fig. 3 caption)
- **Claim:** “The first two principal components explained approximately 98.11% of the variance.”
 - **Checks:** range_check (percentage bounds)
 - **Verdict:** PASS
 - **Notes:** Percentage within [0, 100].
12. ✓ **C12_mass_range_mapping_from_double_log** (p.6, Sec. 3.4 Learned Features)
- **Claim:** “global range for $\log_{10}(\log_{10}(\text{Mass}))$... [0.9829, 1.1619]. This translates to actual halo mass ranging from $10^{9.61}$ to $10^{14.52} (M_{\odot} h^{-1})$.”
 - **Checks:** unit-consistent transform / inverse mapping
 - **Verdict:** PASS
 - **Notes:** Both endpoints match inverse mapping within tolerance.
13. ✓ **C13_actual_mass_endpoints_from_log10_mass** (p.6, Sec. 3.4 Learned Features)
- **Claim:** “actual halo mass ranging from $10^{9.61} M_{\odot} h^{-1}$ to $10^{14.52} M_{\odot} h^{-1}$.”
 - **Checks:** magnitude computation
 - **Verdict:** PASS
 - **Notes:** Ordering and ratio identity hold.
14. ✓ **C14_train_val_test_percent_sum** (p.2, Sec. 2.1 Data Acquisition and Pre-processing)
- **Claim:** “training (80%), validation (10%), and testing (10%) sets.”
 - **Checks:** percent_sum_to_100
 - **Verdict:** PASS
 - **Notes:** Percentages sum to 100.

Limitations

- Checks are limited to arithmetic/logical consistency using only the numeric values explicitly stated in the PDF text; model outputs, dataset statistics, and per-sample metrics cannot be recomputed without the underlying data.
- Figures are not used for extracting numeric values (no pixel/plot reading), so claims that rely on plotted BIC/AIC or other curves cannot be verified here.
- Some statements are qualitative (“significantly lower”, “broad range”) and can only be partially assessed via basic range/magnitude checks rather than full statistical validation.