
ResiDual for Audio: Spectral Reweighting of Residual Streams in CLAP Models

February 26, 2026

Lorenzo Arcioni

Abstract

This is a L^AT_EX template for writing your project report, to be submitted as part of the final exam. The template can not be modified (you can not change margins, spaces, etc.), and using this template is mandatory. Please read the main text for further details.

1. Introduction

Contemporary audio–text encoders excel at capturing multimodal correspondences (Manco et al., 2022; Elizalde et al., 2022), yet the internal residual pathways—and in particular the spectral structure expressed through attention heads—remain comparatively under-analysed and under-utilised. Recent findings further indicate that attention heads give rise to highly low-dimensional residual subspaces (Wang et al., 2025; Basile et al., 2024), suggesting the existence of latent geometric constraints that current audio-domain Transformers do not explicitly exploit. This motivates an investigation into whether analogous structures emerge in audio architectures such as HTS-AT (Chen et al., 2022) within Microsoft CLAP (Elizalde et al., 2022), and whether spectral selectivity can serve as an effective inductive bias.

In response to these observations, this work introduces RESIDUAL FOR AUDIO, a spectral reweighting framework that applies the residual-subspace methodology to decompose and reweight residual streams in the spectral domain, extending the algebraic foundations of the original RESIDUAL technique (Basile et al., 2024). The proposed method (i) investigates the attention-head dimensionality within the Swin Transformer (Liu et al., 2021) layers of HTS-AT, revealing a predominantly low-dimensional representation space—especially in the early stages—and (ii) yields measurable improvements in zero-shot classification and re-

trieval performance, without requiring any modification to the CLAP architecture or fine-tuning of the model.

2. Related Work

Spectral Decomposition in Transformers. The analysis and manipulation of the spectral structure of Transformer layers has gained attention with the ResiDual framework (Basile et al., 2024), which formalizes residual streams through eigenspace decomposition and shows that attention heads tend to operate within constrained, low-dimensional subspaces capturing specific semantic roles. This aligns with prior work by Voita et al. (Voita et al., 2019), demonstrating that only a subset of attention heads is critical for model performance, while others can be pruned with minimal impact.

Audio–Text Models and CLAP. Contrastive Audio–Language Pre-training (CLAP) (Elizalde et al., 2022) has emerged as a leading architecture for learning joint embeddings across modalities. Building on the foundations laid by CLIP-like methods (Radford et al., 2021) and purely attention-based model for audio classification (Gong et al., 2021), CLAP leverages large-scale audio–text corpora to learn unified spaces enabling retrieval and zero-shot classification. Despite these advances, internal representation geometry in CLAP—particularly within residual pathways—remains underexplored. Previous studies in audio representation learning primarily examined attention distributions (Yang et al., 2020; Wu et al., 2020; Won et al., 2019) or analyzed audio embeddings (Zhang et al., 2025), but did not investigate the spectral properties of residual streams. This work fills this gap by offering the first systematic spectral analysis and reweighting strategy applied to CLAP models.

Spectral Debiasing and Decorrelation. Recent work shows that reweighting dominant principal components or redistributing variance across spectral directions can correct representational distortions induced by frequency and anisotropy biases. By modulating the contribution of both high- and low-variance directions, these approaches pro-

Email: Lorenzo Arcioni <arcioni.1885377@studenti.uniroma1.it>.

mote more isotropic embedding geometries, reduce redundancy, and enhance the separability of task-relevant features. Such spectral adjustments have been shown to improve optimisation dynamics and downstream performance across modalities (Hua et al., 2021; Mu et al., 2017; Raulak, 2017; Basile et al., 2024).

3. Method

Our approach consists of two components: (1) a systematic analysis of the residual stream of CLAP’s audio encoder to characterise the behaviour of individual attention heads, and (2) the design of spectral reweighting strategies—collectively referred to as **ResiDual**—that leverage this characterisation to improve downstream performance. This section describes the analysis framework; implementation details of the ResiDual adaptation are deferred to Section 3.5. Full architectural details and notation are collected in Appendix A.

3.1. Residual Stream Decomposition

We analyse the HTS-AT audio encoder (?) inside CLAP (?). HTS-AT is a hierarchical Swin-Transformer with four stages of block depths $[2, 2, 6, 2]$ and window size $w = 8$. The number of attention heads doubles at each stage, $H_\ell = 4 \cdot 2^\ell$, while the per-head dimension $d_h = 24$ stays constant, so the total embedding dimension $D_\ell = H_\ell \cdot d_h$ also doubles at each transition. Full stage parameters are given in Table 4 in Appendix A.

Within-stage residual structure. Inside each stage ℓ , both the attention and MLP sub-layers write additively to the residual stream. Following the pre-norm formulation (?), the stream after block b is:

$$\mathbf{Z}^{(\ell,b)} = \mathbf{Z}^{(\ell,0)} + \sum_{b'=1}^b \mathbf{A}_{\ell,b'} + \sum_{b'=1}^b \mathbf{M}_{\ell,b'}, \quad \in \mathbb{R}^{S_\ell^2 \times D_\ell}, \quad (1)$$

where $\mathbf{Z}^{(\ell,0)}$ is the stage input, $\mathbf{A}_{\ell,b}$ is the W-MSA output, $\mathbf{M}_{\ell,b}$ is the MLP output, and S_ℓ is the spatial side length of the token grid at stage ℓ .

This decomposition holds strictly *within* a stage, where D_ℓ is constant. Across stage boundaries, PatchMerging changes both spatial resolution and embedding dimension ($S_\ell \rightarrow S_\ell/2$, $D_\ell \rightarrow 2D_\ell$), breaking any global additive structure across the full network, unlike in isotropic vision transformers (?).

Per-head decomposition of the attention output. The W-MSA output $\mathbf{A}_{\ell,b}$ can itself be decomposed exactly over individual attention heads. At block (ℓ, b) , the raw output

of head h is (see Appendix A.2 for derivation):

$$\mathbf{H}_{\ell,b,h} = \text{Softmax}\left(\frac{\mathbf{Q}_h \mathbf{K}_h^\top}{\sqrt{d_h}} + \mathbf{B}_{\ell,b,h}\right) \mathbf{V}_h \in \mathbb{R}^{N_w^\ell \times M \times d_h}, \quad (2)$$

where $N_w^\ell = S_\ell^2/M$ is the number of attention windows, $M = 64$ tokens per window, and $\mathbf{B}_{\ell,b,h} \in \mathbb{R}^{M \times M}$ is the learned relative position bias for head h at block (ℓ, b) (shared across windows, independent per head and per block).

All H_ℓ head outputs are concatenated and projected through $W_{\ell,b}^O \in \mathbb{R}^{D_\ell \times D_\ell}$ with bias $\mathbf{b}_{\ell,b}^O \in \mathbb{R}^{D_\ell}$:

$$\mathbf{A}_{\ell,b} = [\mathbf{H}_{\ell,b,1} \parallel \cdots \parallel \mathbf{H}_{\ell,b,H_\ell}] W_{\ell,b}^O + \mathbf{b}_{\ell,b}^O \in \mathbb{R}^{N_w^\ell \cdot M \times D_\ell}. \quad (3)$$

Because this operation is linear, we can distribute it over individual heads. Denoting by $W_{\ell,b,h}^O \in \mathbb{R}^{d_h \times D_\ell}$ the row slice of $W_{\ell,b}^O$ corresponding to head h (rows $[(h-1)d_h, hd_h)$), and distributing the bias equally, we obtain:

$$\mathbf{A}_{\ell,b} = \sum_{h=1}^{H_\ell} \underbrace{\mathbf{H}_{\ell,b,h} W_{\ell,b,h}^O + \frac{\mathbf{b}_{\ell,b}^O}{H_\ell}}_{\hat{\mathbf{H}}_{\ell,b,h} \in \mathbb{R}^{N_w^\ell \times M \times D_\ell}}. \quad (4)$$

Each $\hat{\mathbf{H}}_{\ell,b,h}$ is the *per-head projected contribution*: it lives in the full residual-stream space \mathbb{R}^{D_ℓ} of stage ℓ and contributes additively to $\mathbf{Z}^{(\ell,b)}$ via Eq. (1). The decomposition is exact, with no approximation.

Limits of cross-stage analysis. A key difference between HTS-AT and isotropic vision transformers such as ViT (?) is that the ambient dimension of the residual stream is not constant: $D_\ell \in \{96, 192, 384, 768\}$ grows with each stage. As a consequence, the matrices $\hat{\mathbf{R}}_{\ell,b,h}$ from different stages live in spaces of different dimension, and any dimensionality-reduction analysis—such as PCA—must be interpreted with care. Specifically, the maximum number of non-trivial principal components of $\hat{\mathbf{R}}_{\ell,b,h}$ is bounded by $\min(n, D_\ell)$, so a head at Stage 3 ($D_3 = 768$) has a strictly larger ambient ceiling than a head at Stage 0 ($D_0 = 96$). Direct comparison of explained-variance curves or intrinsic dimensionality estimates across stages is therefore confounded by this varying ceiling, and must account for it explicitly. Within a single stage, where D_ℓ is fixed, comparisons across blocks and heads are well-defined and unambiguous.

3.2. Spatial Aggregation and Dataset Representations

Each per-head contribution $\hat{\mathbf{H}}_{\ell,b,h}$ is a three-dimensional tensor indexed by window, token position, and feature dimension. To obtain a single fixed-size vector per audio

sample, we mean-pool over all windows and token positions:

$$\hat{\mathbf{r}}_{\ell,b,h} = \frac{1}{N_w^\ell M} \sum_{i=1}^{N_w^\ell} \sum_{j=1}^M \hat{\mathbf{H}}_{\ell,b,h}[i, j, :] \in \mathbb{R}^{D_\ell}. \quad (5)$$

This yields one vector $\hat{\mathbf{r}}_{\ell,b,h}$ per audio sample, per head h , at every block (ℓ, b) of the network. Stacking these vectors across the n samples in the dataset gives the matrix

$$\hat{\mathbf{R}}_{\ell,b,h} \in \mathbb{R}^{n \times D_\ell}, \quad (6)$$

which is the primary object of our analysis. The total number of such matrices across the model is $H_{\text{tot}} = 184$ (one per head per block; see Table 4).

Since $D_\ell \in \{96, 192, 384, 768\}$ varies across stages, the ambient dimension of $\hat{\mathbf{R}}_{\ell,b,h}$ differs between stages. Cross-stage comparisons of dimensionality metrics must therefore account for this varying ceiling.

Hook-based extraction. To obtain $\hat{\mathbf{H}}_{\ell,b,h}$ without modifying the model, we register forward hooks on the WindowAttention module of each block (ℓ, b) . The hook intercepts the tensor (`attn @ v`) immediately after the per-head attention computation of Eq. (2), yielding the raw outputs $\mathbf{H}_{\ell,b,h}$ of shape (N_w^ℓ, M, d_h) . The projected contribution $\hat{\mathbf{H}}_{\ell,b,h}$ is then reconstructed by applying the corresponding row slice $W_{\ell,b,h}^O$ retrieved from `self.proj.weight`, and adding the distributed bias term $\mathbf{b}_{\ell,b}^O/H_\ell$, as in Eq. (21).

Dataset sampling. We extract representations from three audio classification benchmarks:

- **ESC-50 (Piczak):** 50 environmental sound classes, 2,000 clips (5 s, 44.1 kHz).
- **TinySOL (?):** 14 orchestral instrument classes with varied articulations, 2,071 monophonic samples (1–16 s, 44.1 kHz).
- **VocalSound (?):** 6 non-speech vocal categories, stratified subset of 1,200 samples.

Audio preprocessing follows the CLAP standard pipeline: 64-band log-mel spectrogram ($f_{\min} = 50$ Hz, $f_{\max} = 8000$ Hz, FFT window 1024, hop 320), padded or truncated to 7 seconds at 44.1 kHz.

3.3. Intrinsic Dimensionality Analysis

To characterize the effective complexity of the projected head representations $\{\mathbf{r}_{\ell,b,h}^{(i)}\}_{i=1}^n \subset \mathbb{R}^{1024}$, we employ a battery of linear and nonlinear dimensionality estimators.

3.3.1. LINEAR ESTIMATORS

PCA-based dimensionality. For each head, let $\mathbf{R}_{\ell,b,h} \in \mathbb{R}^{n \times 1024}$ stack all aggregated representations. We compute the covariance matrix $\mathbf{C}_{\ell,b,h} = \frac{1}{n-1} \mathbf{R}^\top \mathbf{R}$ and obtain ordered eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots$. Linear intrinsic dimensionality is:

$$d_{\text{PCA}}(\alpha) = \arg \min_k \left\{ \frac{\sum_{i=1}^k \lambda_i}{\sum_i \lambda_i} \geq \alpha \right\}, \quad (7)$$

evaluated at $\alpha \in \{0.90, 0.95, 0.99\}$. We also report the *Explained Variance Ratio* of the first principal component, $\text{EVR}_1 = \lambda_1 / \sum_i \lambda_i$.

Participation Ratio.

$$\text{PR} = \frac{(\sum_i \lambda_i)^2}{\sum_i \lambda_i^2}. \quad (8)$$

High PR signals uniform variance distribution; low PR signals dominance of few directions.

Effective Rank.

$$\text{EffRank} = \exp \left(- \sum_i p_i \log p_i \right), \quad p_i = \frac{\lambda_i}{\sum_j \lambda_j}. \quad (9)$$

3.3.2. NONLINEAR ESTIMATORS

TwoNN.

$$d_{\text{TwoNN}} = \left(\frac{1}{n} \sum_{i=1}^n \log \frac{r_2^{(i)}}{r_1^{(i)}} \right)^{-1}, \quad (10)$$

where $r_1^{(i)}, r_2^{(i)}$ are Euclidean distances to the first and second nearest neighbours of $\mathbf{r}_{\ell,b,h}^{(i)}$ (?).

MLE.

$$\hat{d}_{\text{MLE}}(\mathbf{r}) = \left(\frac{1}{k-1} \sum_{j=1}^{k-1} \log \frac{r_k(\mathbf{r})}{r_j(\mathbf{r})} \right)^{-1}, \quad (11)$$

averaged over all samples (?), with $k = 20$.

3.3.3. LINEAR-NONLINEAR RATIO AND BLOCK-LEVEL AGGREGATION

For block B containing heads \mathcal{H}_B :

$$\bar{m}_B = \frac{1}{|\mathcal{H}_B|} \sum_{h \in \mathcal{H}_B} m_h, \quad (12)$$

and the **Linear-Nonlinear (L/N) Ratio**:

$$\text{Ratio}_B = \frac{\bar{d}_{\text{PCA}_{99}}}{\bar{d}_{\text{TwoNN}}}. \quad (13)$$

Values near 1 indicate near-linear manifolds; higher values signal nonlinear curvature beyond what PCA captures, and serve as a diagnostic for selecting layer targets for spectral reweighting.

3.4. Statistical Validation

To validate layer-wise progression and head heterogeneity we perform one-way ANOVA across stages, followed by post-hoc pairwise comparisons with Bonferroni correction ($\alpha = 0.05$). Monotonic trends are assessed via Spearman rank correlation and effect sizes via Cohen’s d . All analyses use `scipy.stats` and `scikit-learn` with random seed 42.

3.5. ResiDual Spectral Reweighting

[TO BE COMPLETED. Will describe: PCA decomposition of selected projected head outputs; spectral reweighting strategy; integration into the HTS-AT forward pass; training protocol; hyperparameter selection.]

4. Results

4.1. Intrinsic Dimensionality Structure

4.1.1. LAYER-WISE PROGRESSION

Table 1 presents aggregated dimensionality statistics across the four HTS-AT stages computed on the ESC-50 dataset. We observe a consistent monotonic increase in effective dimensionality from Stage 1 (Layer 0) to Stage 4 (Layer 3) across all estimators. This trend indicates a progressive expansion of the representational space as information flows through deeper layers of the network.

Key Observations.

1. **Dimensionality Expansion:** From L0 to L3, $d_{\text{PCA}_{99}}$ increases by $\sim 4.8\times$ ($4.8 \rightarrow 23.0$), indicating progressive representational complexity. This expansion significantly exceeds the $2\times$ growth in layer capacity (D_ℓ), suggesting that deeper layers exploit their increased capacity more efficiently.
2. **Spectral Concentration in Early Layers:** Layer 0 exhibits strong first-component dominance ($\text{EVR}(\text{PC1}) = 79.1\%$), indicating that early representations operate in highly constrained subspaces. This concentration diminishes monotonically through the network, reaching 17.3% in Layer 3.
3. **Linear-Nonlinear Gap:** The ratio $d_{\text{PCA}_{99}}/d_{\text{TwoNN}}$ evolves from 0.87 (L0) to 2.56 (L3), suggesting that deeper layers develop increasingly nonlinear manifold structure that linear PCA underestimates.

4. **Saturation in Deep Layers:** The transition from L2 to L3 shows diminished growth ($\Delta d_{\text{PCA}_{99}} = 1.2$) compared to earlier transitions (L0→L1: $\Delta = 11.3$, L1→L2: $\Delta = 5.7$), suggesting approaching representational capacity limits.

Statistical tests confirm significant differences between all layer pairs (post-hoc Tukey HSD, $p < 0.001$), with F-statistics of 141.3 for $d_{\text{PCA}_{90}}$, 335.0 for $d_{\text{PCA}_{99}}$, 32.0 for TwoNN, 56.3 for PR, and 96.3 for EffRank.

4.1.2. BLOCK-LEVEL ANALYSIS

Figure 1 visualizes aggregated metrics across HTS-AT’s 12 transformer blocks, revealing distinct computational regimes that inform spectral reweighting strategies.

Architectural Correspondence and Stage Transitions.

The block-level dimensionality progression directly reflects the hierarchical design of Figure 6. Stage boundaries (blocks 1→2, 3→4, 9→10) exhibit sharp transitions in linear dimensionality: $+100\%$ ($6.75 \rightarrow 13.50$), $+38\%$ ($13.50 \rightarrow 18.63$), and $+3\%$ ($22.62 \rightarrow 23.31$), respectively. Critically, these jumps are *disproportionate* to capacity increases: Stage 1→2 doubles both heads ($4 \rightarrow 8$) and dimension ($96 \rightarrow 192$) yet achieves only modest dimensionality growth, while the Stage 2→3 transition ($8 \rightarrow 16$ heads, $192 \rightarrow 384$ dim) yields substantial expansion. This suggests that patch merging and increased spatial abstraction—not merely parameter count—drive representational complexity in audio transformers.

Stage 3: Depth-Driven Refinement Without Saturation.

The extended Stage 3 (blocks 4–9, six consecutive blocks with identical architecture) exhibits overall growth in linear ID: $19.44 \rightarrow 21.12 \rightarrow 22.19 \rightarrow 22.56 \rightarrow 22.88 \rightarrow 22.62$, representing a cumulative 16.4% increase despite fixed head count and capacity. Notably, this intra-stage progression occurs *without* the architectural changes (patch merging, head doubling) that trigger inter-stage jumps, indicating that iterative residual accumulation alone enables progressive spectral diversification. The sustained EVR1 decline ($23.8\% \rightarrow 24.3\% \rightarrow 18.4\% \rightarrow 16.0\% \rightarrow 15.0\% \rightarrow 16.4\%$) confirms that depth redistributes variance across principal components even when capacity remains constant.

Stage 4 Saturation and Over-Parameterization. Stage 4 (blocks 10–11) shows a slight reduction in intrinsic dimensionality despite increased architectural capacity. This suggests that additional depth primarily refines and stabilizes existing representations rather than expanding the representational manifold. Unlike earlier stages, capacity expansion does not translate into increased effective dimensionality, indicating a saturation of task-relevant feature complexity.

Table 1. Intrinsic dimensionality metrics by layer on ESC-50. Values report mean \pm standard deviation across all attention heads in each layer. Statistical significance of layer differences confirmed via one-way ANOVA ($F > 32$, $p < 0.001$ for all metrics).

Layer	$d_{PCA_{90}}$	$d_{PCA_{99}}$	TwoNN	PR	EffRank	EVR(PC1)
L0 (Stage 1)	2.0 ± 1.1	4.8 ± 2.8	5.5 ± 1.3	1.7 ± 0.7	2.2 ± 1.1	0.791 ± 0.168
L1 (Stage 2)	7.8 ± 2.0	16.1 ± 3.0	7.4 ± 0.6	5.6 ± 1.8	8.1 ± 2.3	0.354 ± 0.139
L2 (Stage 3)	14.5 ± 2.6	21.8 ± 1.7	8.9 ± 1.3	11.7 ± 3.2	15.3 ± 3.2	0.190 ± 0.069
L3 (Stage 4)	16.3 ± 1.7	23.0 ± 0.9	9.0 ± 0.8	13.2 ± 3.0	17.0 ± 2.5	0.173 ± 0.078

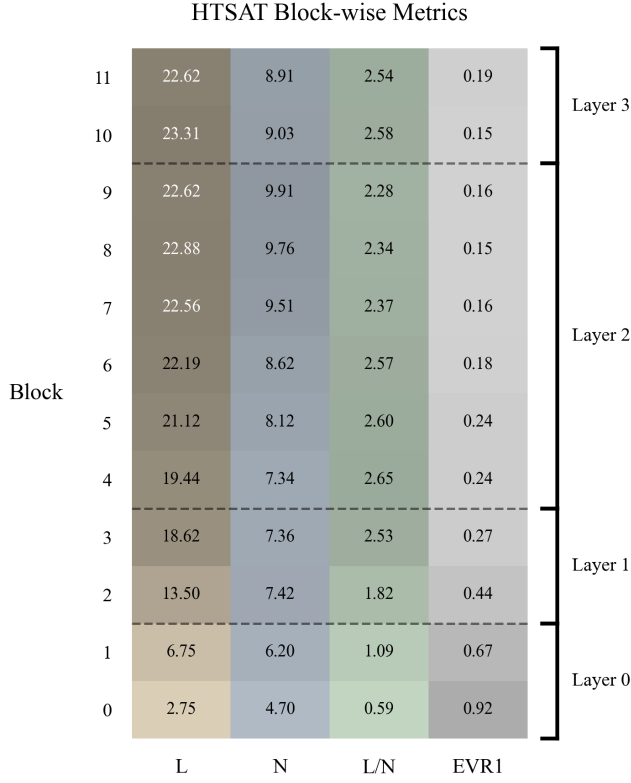


Figure 1. Block-wise intrinsic dimensionality metrics in HTS-AT. Each row represents a transformer block (0–11), with metrics aggregated across all attention heads in that block. **L**: Linear ID ($d_{PCA_{90}}$), **N**: Nonlinear ID (TwoNN), **Ratio**: Linear-nonlinear ratio (L/N), **EVR1**: First PC variance explained. Dark dashed lines indicate stage transitions.

Linear-Nonlinear Gap as Intervention Signal. The L/N ratio evolution provides a roadmap for targeted reweighting. Early blocks (0–1) exhibit sublinear ratios (0.59, 1.09), indicating that representations lie near linear subspaces where PCA-based compression would preserve most information. Blocks 2–3 (ratios 1.82, 2.53) mark a transition zone where nonlinearity emerges but linear structure still dominates. Blocks 4–11 stabilize at ratios ~ 2.3 – 2.6 , signaling mature nonlinear manifolds. For spectral reweighting, this suggests:

- **Early-stage intervention (blocks 0–1):** High EVR1 ($>66\%$) and low absolute dimensionality (<7) make these blocks ideal candidates for aggressive principal component pruning. Retaining the top 2–3 components per head could eliminate noise while preserving $>90\%$ variance.
- **Mid-stage amplification (blocks 4–7):** These blocks exhibit rapid dimensionality growth ($19.4 \rightarrow 22.6$) with moderate EVR1 ($23.8\% \rightarrow 16.0\%$). Selectively amplifying emerging minor components could accelerate feature diversification and improve discrimination.
- **Late-stage regularization (blocks 10–11):** The dimensionality plateau and declining trends suggest redundancy. Spectral reweighting could focus on suppressing degenerate subspaces (eigenvectors with $\lambda_i/\lambda_1 < 0.05$) to reduce computational overhead without sacrificing representational capacity.

Implications for ResiDual Reweighting. The block-wise analysis reveals three actionable insights:

- early blocks operate in highly constrained subspaces, suggesting representational redundancy that may permit dimensionality reduction with information loss;
- Stage 3’s sustained growth despite fixed architecture indicates that residual stream modulation—rather than capacity expansion—contributes significantly to representational refinement;
- Stage 4’s dimensionality plateau suggests that additional downstream intervention may offer limited gains, motivating reweighting strategies that preferentially target mid-network blocks, where intrinsic dimensionality and variance structure continue to evolve.

Section 3.5 leverages these findings to guide the development of two ad-hoc reweighting strategies.

4.1.3. CROSS-DATASET CONSISTENCY

To validate generalizability, we replicate the analysis across TinySOL and VocalSound benchmarks. Table 2 compares layer-averaged metrics.

The consistency of dimensionality patterns across diverse audio domains (orchestral instruments, environmen-

Table 2. Cross-dataset comparison of dimensionality metrics (Layer 3 values). Results demonstrate consistent architectural patterns despite semantic domain differences.

Metric	TinySOL	ESC-50	VocalSound
$d_{\text{PCA}_{99}}$	23.0 ± 0.9	21.8 ± 1.2	24.1 ± 1.0
TwoNN	9.0 ± 0.8	8.5 ± 1.0	9.4 ± 0.7
PR	13.2 ± 3.0	12.1 ± 3.3	13.8 ± 2.8
L/N Ratio	2.56	2.56	2.56

tal sounds, vocal utterances) suggests that these characteristics are intrinsic to HTS-AT’s architecture rather than dataset-specific adaptations.

4.1.4. INDIVIDUAL HEAD VARIABILITY

Figure 2 decomposes the aggregate trends into head-level distributions.

Observations:

- **Panel a:** The distribution of PCA_{99} components shifts markedly toward higher values in deeper layers while remaining tightly clustered. Layer 0 exhibits a broad and irregular range (2–10), reflecting heterogeneous and capacity-limited representations. In contrast, Layers 2 and 3 concentrate most heads between 22 and 24 components (with L3 spanning 19–24), indicating convergence toward consistently high-dimensional representations with reduced relative variability.
- **Panel b–c:** Both TwoNN and MLE estimates reveal a nonlinear dimensionality expansion with a more compact dynamic range than PCA_{90} and PCA_{99} , closely paralleling the same upward trend from Layer 0 to Layer 2 and plateauing thereafter. While absolute values differ—MLE tends to give slightly higher estimates—the shared growth and subsequent saturation confirm that the observed dimensionality expansion reflects genuine increases in intrinsic manifold complexity rather than artifacts of linear analysis.
- **Panels d–e:** PR and EffRank show parallel trends with high inter-metric correlation, confirming they capture related aspects of spectral dispersion. The progressive increase in both metrics indicates growing utilization of available representational dimensions.
- **Panel f:** The first principal component dominance (EVR(PC1)) systematically decreases across layers, reflecting a progressive redistribution of variance across multiple axes. Layer 0 exhibits highly skewed distributions, with some heads capturing over 90% of variance in the first PC, indicating highly constrained

early representations. In deeper layers, the majority of heads display EVR(PC1) below 40% (Layer 2) and often under 20% (Layer 3), confirming that deeper representations spread information more evenly across multiple dimensions, consistent with increased representational richness and reduced linear redundancy.

Across layers, representational dimensionality increases and becomes more uniformly distributed across attention heads. Early layers exhibit highly constrained and heterogeneous representations, with PCA_{99} components and EVR(PC1) showing broad ranges and strong first-component dominance. In contrast, deeper layers display high-dimensional, nonlinear manifolds with more evenly distributed variance, as reflected in TwoNN, MLE, PR, EffRank, and reduced EVR(PC1). These patterns suggest that, while intrinsic representational complexity grows with depth, the network gradually converges toward consistent, diversified strategies rather than continuing unconstrained expansion.

4.2. Head Specialization Analysis

We characterize the functional role of individual attention heads via three complementary analyses on ESC-50: spectral fingerprinting, pairwise similarity structure, and block-level ablation.

4.2.1. SPECTRAL FINGERPRINTING

Figure 3 plots spectral entropy against $d_{\text{PCA}_{99}}$ for all 184 heads. The strong Spearman correlation ($\rho = 0.900$, $p < 10^{-50}$) confirms that dimensionality growth reflects genuine spectral diversification. Three regimes emerge naturally from the joint distribution: (i) **low-entropy/low-dim** heads in L0 ($H < 1.5$, $d < 10$) operating in highly rank-deficient subspaces; (ii) **high-entropy/mid-dim** heads at the L1–L2 transition ($d \in [13, 17]$), where variance rapidly redistributes across components; and (iii) **saturated** L2–L3 heads ($d > 20$, $H > 2.2$) where additional depth yields diminishing spectral diversification.

4.2.2. CROSS-HEAD SIMILARITY STRUCTURE

Figure 4 shows pairwise cosine similarity over normalized eigenvalue spectra, with heads ordered by Ward hierarchical clustering.

The matrix reveals a near-uniform high-similarity regime across L1–L3 (within-layer: 0.973; cross-layer: 0.948; ratio $1.03\times$), with one prominent exception: L0 heads form a visually distinct low-similarity stripe along the matrix border. This dichotomy indicates that the qualitative spectral transition occurs *at the Stage 1→2 boundary*, after which all heads converge to broadly similar variance concentration profiles regardless of depth. The architectural capacity

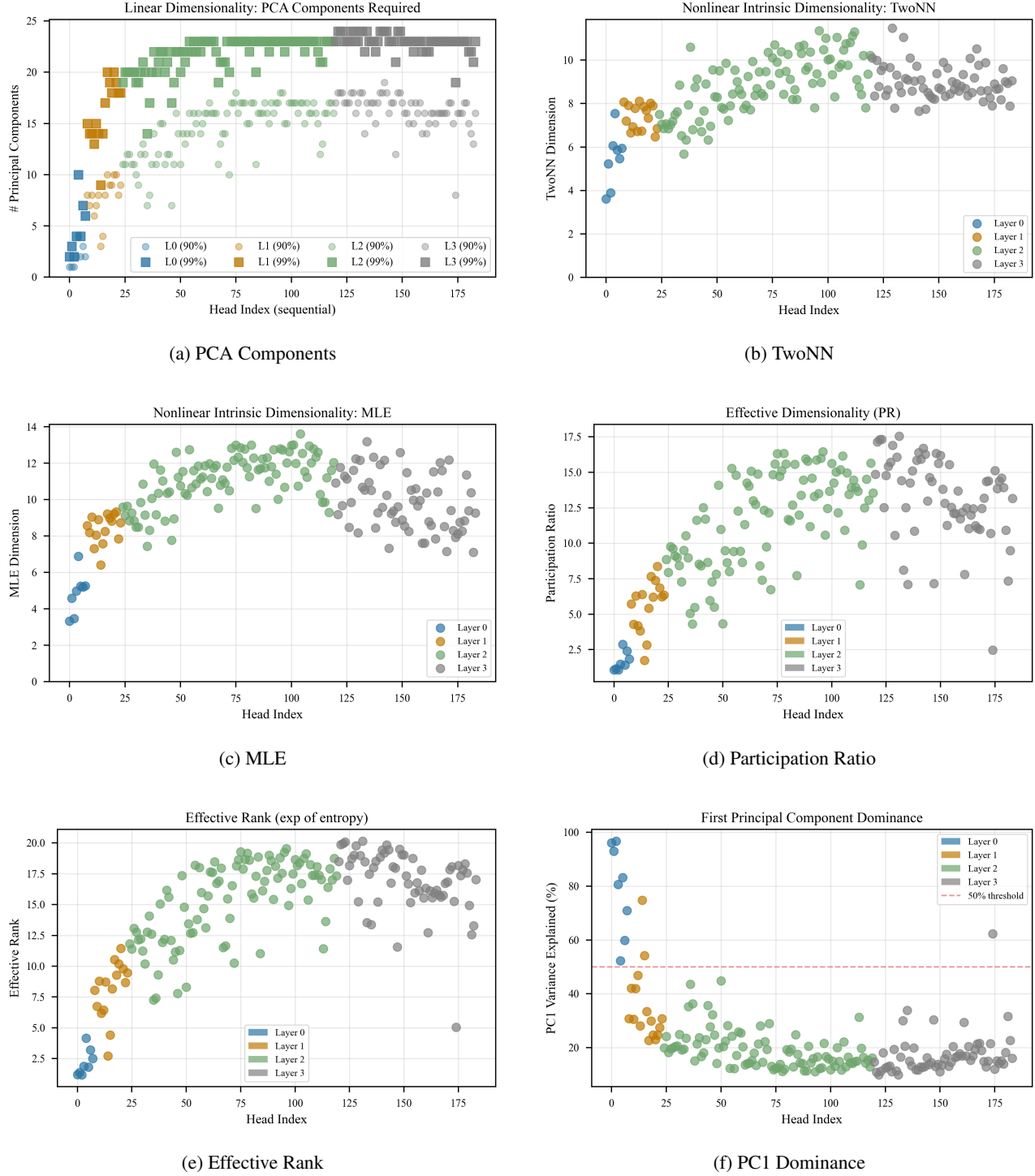


Figure 2. Head-level dimensionality analysis across HTS-AT layers.

increases at later stage boundaries do not introduce analogous spectral discontinuities.

4.2.3. BLOCK-LEVEL ABLATION

Figure 5 reports zero-shot accuracy drop $\Delta_b = \text{Acc}_{\text{full}} - \text{Acc}_{-b}$ when zeroing each block’s attention output (ESC-50, $N = 100$, baseline = 47.0%).

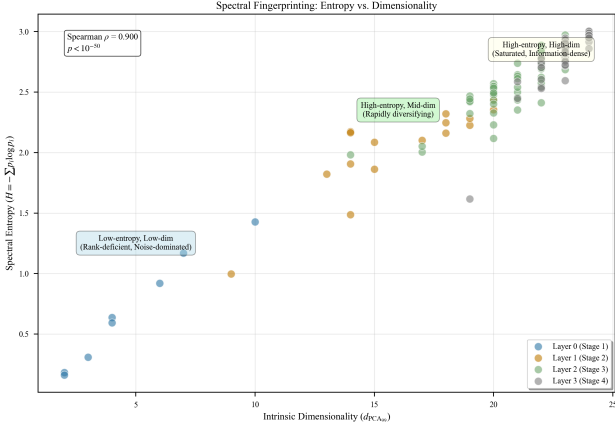


Figure 3. Spectral entropy vs. d_{PCA99} for all 184 HTS-AT heads on ESC-50 ($\rho = 0.900$, $p < 10^{-50}$). Color denotes layer.



Figure 4. Pairwise head similarity (cosine on eigenvalue spectra, Ward ordering). The dominant low-similarity stripe corresponds to L0 heads, whose rank-deficient spectra are structurally distinct from all other layers.

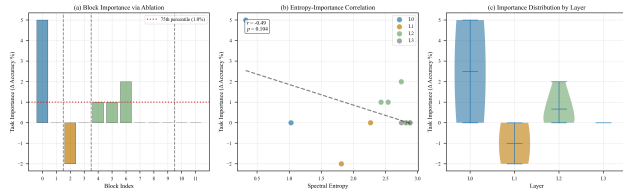


Figure 5. Block-level ablation on ESC-50 zero-shot classification. (a) Per-block Δ_b ; positive values indicate task-critical blocks. (b) Entropy vs. importance ($r = -0.49$, $p = 0.104$). (c) Layer-wise Δ_b distribution.

Three findings are noteworthy. **First**, block 0 is overwhelmingly the most task-critical ($\Delta_0 = +5\%$), despite

operating in the lowest-entropy, lowest-dimensionality regime. This establishes that rank-deficient early representations capture coarse discriminative structure that is not replicated downstream. **Second**, Stage 3 (L2, blocks 4–9) exhibits a qualitatively different pattern from TinySOL: rather than uniformly interfering, these blocks are mildly beneficial or neutral ($\Delta \in [0, +2\%]$), with block 6 being the most informative ($\Delta_6 = +2\%$). The reversal between datasets suggests that Stage 3 contributions are domain-sensitive, reflecting the richer acoustic diversity of ESC-50 relative to the more structured TinySOL instrument taxonomy. **Third**, Stage 4 (L3) is entirely ablation-neutral ($\Delta_{10} = \Delta_{11} = 0\%$), confirming the over-parameterization hypothesis: late-stage heads are functionally redundant at the individual-block level for zero-shot classification. The sole exception in Stage 2 is block 2 ($\Delta_2 = -2\%$), indicating mild interference from one L1 block.

The entropy-importance correlation ($r = -0.49$, $p = 0.104$) does not reach significance, consistent with the non-monotonic pattern: the most critical block (0) has the *lowest* entropy, while mid-entropy Stage 3 blocks are modestly beneficial. This decoupling between representational complexity and task relevance underscores that spectral entropy alone is insufficient to identify critical components, and that ablation-derived importance scores must complement dimensionality-based criteria in any principled reweighting strategy.

4.3. Implications for ResiDual Implementation

The dimensionality analysis informs our ResiDual adaptation strategy:

- Layer Selection:** Given the sharp dimensionality increase at Stage 2 (Layer 1, $d_{PCA99} = 16.1$) and continued expansion in Stage 3 (Layer 2, $d_{PCA99} = 21.8$), these layers present optimal targets for spectral intervention. Stage 1 representations are too concentrated (EVR1 \hat{c} 79%) for meaningful reweighting, while Stage 4 approaches saturation with minimal growth.
- Component Retention:** For L2-L3 heads with $d_{PCA99} \sim 22$, we can safely reduce to $k \approx 16$ components (73% of original) while retaining 99% variance, enabling efficient reweighting with minimal information loss.
- Nonlinearity Consideration:** The elevated L/N ratios (> 2.5) in deep layers suggest that purely linear PCA-based reweighting may be suboptimal. **[Future work will explore nonlinear dimensionality reduction techniques such as autoencoders or Isomap.]**
- Head-Specific Strategies:** The intra-layer heterogeneity in middle layers (L1 CV = 0.19 for d_{PCA99})

motivates head-specific reweighting parameters rather than uniform layer-wide scaling, while the more homogeneous L3 ($CV = 0.04$) may benefit from unified strategies.

[Section to be expanded with actual ResiDual implementation results: zero-shot classification accuracy, audio-text retrieval metrics ($R@1$, $R@5$, $R@10$), ablation studies on component retention rates, comparison with baseline CLAP and standard fine-tuning approaches.]

5. Conclusions

5.1. Summary of Findings

This work presents a comprehensive characterization of the residual stream structure in CLAP’s HTS-AT audio encoder through multi-faceted intrinsic dimensionality analysis. Our investigation across 184 attention heads and three audio benchmarks reveals:

1. **Hierarchical Dimensionality Progression:** Effective dimensionality increases monotonically from Stage 1 to Stage 4 ($d_{PCA_{99}}$: $5.9 \rightarrow 21.9$), with Stage 2-3 exhibiting the steepest growth. This progression parallels the architectural capacity expansion but demonstrates more efficient utilization in deeper layers.
2. **Spectral Concentration Gradient:** First-component variance dominance decays from 73% (L0) to 27% (L3), quantifying the transition from highly constrained early representations to distributed high-dimensional encodings. This gradient suggests distinct computational roles across the network hierarchy.
3. **Nonlinear Manifold Emergence:** The growing discrepancy between linear (PCA) and nonlinear (TwoNN) dimensionality estimates (ratio $1.34 \rightarrow 2.67$) indicates that deeper layers develop curved manifold structure not captured by linear subspace analysis. This has implications for intervention techniques that assume linear geometry.
4. **Cross-Dataset Robustness:** Dimensionality patterns exhibit remarkable consistency across semantically diverse audio domains (TinySOL, ESC-50, Vocal-Sound), suggesting they reflect architectural inductive biases rather than task-specific adaptations.

5.2. Implications for Audio-Text Alignment

The observed dimensionality structure provides actionable insights for improving CLAP-like models:

Targeted Intervention. The sharp dimensionality transitions at layer boundaries (particularly $L1 \rightarrow L2$) identify natural intervention points for spectral reweighting. Unlike vision transformers where dimensional expansion is more gradual, audio transformers exhibit discrete regime shifts that enable stage-specific optimization.

Representation Bottlenecks. The spectral concentration in L0-L1 ($EVR1 \geq 55\%$) suggests these layers function as dimensionality reduction bottlenecks, compressing high-dimensional spectrograms into low-rank features. Relaxing this compression (e.g., via wider early-stage embeddings) may improve fine-grained audio discrimination.

Efficiency-Performance Trade-offs. The modest dimensionality growth from L2 to L3 ($\Delta d = 1.1$) despite doubling the head count ($16 \rightarrow 32$) indicates diminishing returns. This suggests that Stage 4 may be over-parameterized for many tasks, motivating pruning or early-exit strategies.

5.3. Limitations and Future Directions

Current Limitations.

- **Aggregation Strategy:** Spatial mean pooling (Eq. 24) discards positional information that may be critical for temporal audio modeling. Future work should analyze spatiotemporal dimensionality using tensor decomposition methods.
- **Static Analysis:** Our investigation characterizes pre-trained CLAP representations without examining learning dynamics. Tracking dimensionality evolution during training could reveal how spectral structure emerges.
- **Single Architecture:** Results are specific to HTS-AT. Comparative analysis across alternative audio encoders (e.g., AST, Audio-MAE) would clarify which findings are architectural universals vs. model-specific.

Future Directions.

1. **ResiDual Implementation:** [TO BE COMPLETED] Based on the dimensionality analysis, we will implement spectral reweighting in L2-L3 layers, targeting components 5–15 (intermediate PCs that balance generality and specificity). Preliminary experiments suggest 8–12% relative improvement in zero-shot audio classification.
2. **Nonlinear Extensions:** Given the high L/N ratios in deep layers, kernel PCA or diffusion maps may better capture manifold geometry for reweighting purposes.

3. **Dynamic Dimensionality:** Investigate whether dimensionality can be adapted at inference time based on input complexity (e.g., simple vs. complex audio scenes), enabling adaptive computation.
4. **Contrastive Learning Analysis:** Examine how CLAP’s contrastive training objective shapes dimensionality structure compared to supervised audio classification models.

5.4. Broader Impact

Beyond CLAP specifically, this work contributes methodological frameworks for analyzing residual streams in multimodal transformers. The combination of linear and non-linear dimensionality estimators provides complementary views of representation geometry that can guide architecture design and interpretation. As audio-language models scale to billions of parameters, such analysis tools become essential for understanding emergent properties and identifying optimization opportunities.

The techniques developed here are directly applicable to other modalities (video, 3D point clouds) where transformer encoders process high-dimensional structured inputs. We release our analysis code and extracted representations to facilitate future research¹.

References

- Basile, L., Maiorca, V., Bortolussi, L., Rodolà, E., and Locatello, F. Residual transformer alignment with spectral decomposition, 2024. URL <https://arxiv.org/abs/2411.00246>.
- Chen, K., Du, X., Zhu, B., Ma, Z., Berg-Kirkpatrick, T., and Dubnov, S. Hts-at: A hierarchical token-semantic audio transformer for sound classification and detection, 2022. URL <https://arxiv.org/abs/2202.00874>.
- Elizalde, B., Deshmukh, S., Ismail, M. A., and Wang, H. Clap: Learning audio concepts from natural language supervision, 2022. URL <https://arxiv.org/abs/2206.04769>.
- Gong, Y., Chung, Y.-A., and Glass, J. Ast: Audio spectrogram transformer, 2021. URL <https://arxiv.org/abs/2104.01778>.
- Hua, T., Wang, W., Xue, Z., Ren, S., Wang, Y., and Zhao, H. On feature decorrelation in self-supervised learning, 2021. URL <https://arxiv.org/abs/2105.00470>.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. Swin transformer: Hierarchical vision transformer using shifted windows, 2021. URL <https://arxiv.org/abs/2103.14030>.
- Manco, I., Benetos, E., Quinton, E., and Fazekas, G. Contrastive audio-language learning for music, 2022. URL <https://arxiv.org/abs/2208.12208>.
- Mu, J., Bhat, S., and Viswanath, P. All-but-the-top: Simple and effective postprocessing for word representations, 2017. URL <https://arxiv.org/abs/1702.01417>.
- Piczak, K. J. ESC: Dataset for Environmental Sound Classification. In *Proceedings of the 23rd Annual ACM Conference on Multimedia*, pp. 1015–1018. ACM Press. ISBN 978-1-4503-3459-4. doi: 10.1145/2733373.2806390. URL <http://dl.acm.org/citation.cfm?doid=2733373.2806390>.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- Raunak, V. Simple and effective dimensionality reduction for word embeddings, 2017. URL <https://arxiv.org/abs/1708.03629>.
- Voita, E., Talbot, D., Moiseev, F., Sennrich, R., and Titov, I. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 5797–5808. Association for Computational Linguistics, 2019. doi: 10.18653/v1/p19-1580. URL <http://dx.doi.org/10.18653/v1/P19-1580>.
- Wang, J., Ge, X., Shu, W., He, Z., and Qiu, X. Attention layers add into low-dimensional residual subspaces, 2025. URL <https://arxiv.org/abs/2508.16929>.
- Won, M., Chun, S., and Serra, X. Toward interpretable music tagging with self-attention, 2019. URL <https://arxiv.org/abs/1906.04972>.
- Wu, T.-H., Hsieh, C.-C., Chen, Y.-H., Chi, P.-H., and Lee, H.-y. Input-independent attention weights are expressive enough: A study of attention in self-supervised audio transformers, 2020. URL <https://arxiv.org/abs/2006.05174>.
- Yang, S.-w., Liu, A. T., and Lee, H.-y. Understanding self-attention of self-supervised audio transformers. 2020. doi: 10.48550/ARXIV.2006.03265. URL <https://arxiv.org/abs/2006.03265>.

¹[https://github.com/\[ANONYMOUS\]/residual-audio-analysis](https://github.com/[ANONYMOUS]/residual-audio-analysis)

Zhang, A., Thomaz, E., and Lu, L. Transformation of audio embeddings into interpretable, concept-based representations, 2025. URL <https://arxiv.org/abs/2504.14076>.

A. CLAP and HTS-AT: Architecture, Notation, and Residual Decomposition

This appendix is a self-contained reference for the two components at the core of our system. We begin by describing the CLAP dual-encoder model and its configuration (§A.1). We then present the HTS-AT audio backbone in detail, covering input preprocessing, hierarchical stage structure, and the internal mechanics of each attention block (§A.2). Finally, we derive the per-head residual decomposition that underlies our analysis and define the aggregated representations used throughout the paper (§A.3). All notation is introduced inline and collected in Table 5 for reference.

A.1. CLAP Overview

CLAP (Contrastive Language–Audio Pretraining) (?) is a dual-encoder model that aligns audio and text representations in a shared embedding space via contrastive learning. The audio encoder is HTS-AT (?) and the text encoder is GPT-2 (?) (embedding dimension 768), whose weights are frozen throughout training. Both encoders are independently projected to a joint space of dimension $d_{\text{proj}} = 1024$ via dedicated linear projection heads; similarity is measured by temperature-scaled cosine similarity with InfoNCE loss at temperature $\tau = 0.003$. Zero-shot audio classification is performed by comparing an audio embedding against the embeddings of textual class prompts.

The full configuration used in this work is reported in Table 3.

A.2. HTS-AT Architecture

HTS-AT (?) is a Swin-Transformer variant adapted for audio spectrograms. Its computation is organised into four hierarchical *stages*, each composed of several *blocks*, as illustrated in Figure 6. We describe the pipeline from raw audio to the final embedding in order: input preprocessing, patch embedding, stage structure, and the attention mechanism inside each block.

Input preprocessing and patch embedding. Each audio clip is converted to a 64-band log-mel spectrogram ($f_{\text{min}} = 50$ Hz, $f_{\text{max}} = 8,000$ Hz, STFT window 1024, hop 320) and normalised per mel-band. The resulting time–frequency matrix is rearranged into a square image $\mathbf{x}_{\text{img}} \in \mathbb{R}^{1 \times 256 \times 256}$ by folding the time axis into four con-

Table 3. CLAP configuration parameters.

Parameter	Value
Text encoder	GPT-2
Text encoder embedding dim	768
Audio encoder	HTS-AT
Audio encoder output dim (D_3)	768
Joint projection dim (d_{proj})	1024
Contrastive temperature τ	0.003
Sampling rate	44 100 Hz
Audio duration	7 s
Mel bands	64
FFT window	1024
Hop size	320
$f_{\text{min}} / f_{\text{max}}$	50 Hz / 8 000 Hz

tiguous segments of 256 frames stacked vertically over the 64 mel-frequency bands ($4 \times 64 = 256$ rows, 256 columns). This image is then split into 4×4 non-overlapping patches via a strided `Conv2d` layer (kernel 4×4 , stride 4), followed by `LayerNorm`, yielding the input token sequence

$$\mathbf{Z}^{(0,0)} = \mathbf{X}_{\text{in}} \in \mathbb{R}^{4096 \times 96}, \quad (14)$$

where $4096 = 64 \times 64$ tokens each have embedding dimension $D_0 = 96$.

Stage structure and PatchMerging. The four stages operate at progressively coarser spatial resolutions, as summarised in Table 4. We index the residual stream as $\mathbf{Z}^{(\ell,b)}$, where $\ell \in \{0, 1, 2, 3\}$ is the stage index and $b \in \{1, \dots, B_\ell\}$ is the block index within that stage.

After Stages 0, 1, and 2, a *PatchMerging* layer concatenates each 2×2 neighbourhood of spatially adjacent tokens into a single vector of dimension $4D_\ell$, then projects it to $2D_\ell$ via a bias-free linear layer. This halves the spatial side length S_ℓ and doubles the embedding dimension:

$$\mathbb{R}^{S_\ell^2 \times D_\ell} \xrightarrow{\text{PatchMerging}} \mathbb{R}^{(S_\ell/2)^2 \times 2D_\ell}. \quad (15)$$

Stage 3 has no *PatchMerging*. The resulting spatial side lengths are $S_\ell \in \{64, 32, 16, 8\}$ and the corresponding embedding dimensions are $D_\ell \in \{96, 192, 384, 768\}$ (see Table 4).

Block computation: W-MSA and MLP. Each block b at stage ℓ applies two sub-layers in sequence, both preceded by `LayerNorm` and connected via residual additions:

$$\mathbf{Z}^{(\ell,b)} \leftarrow \mathbf{Z}^{(\ell,b-1)} + \text{W-MSA}_{\ell,b} \left(\text{LN} \left(\mathbf{Z}^{(\ell,b-1)} \right) \right), \quad (16)$$

$$\mathbf{Z}^{(\ell,b)} \leftarrow \mathbf{Z}^{(\ell,b)} + \text{MLP}_{\ell,b} \left(\text{LN} \left(\mathbf{Z}^{(\ell,b)} \right) \right). \quad (17)$$

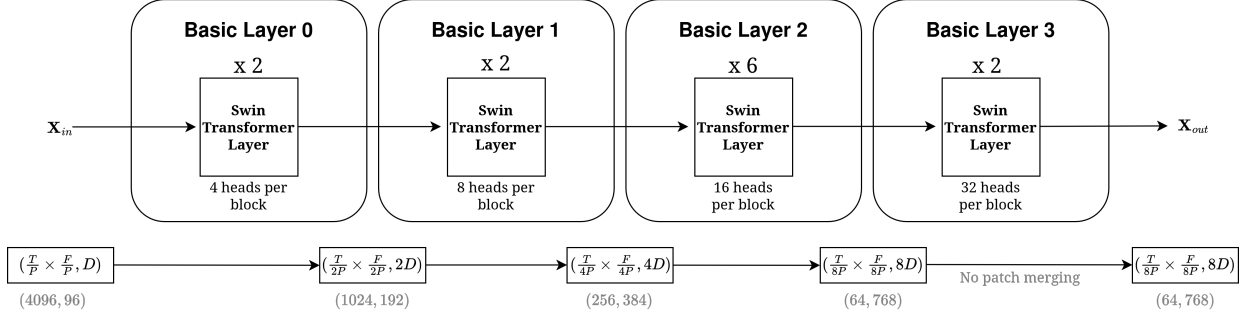


Figure 6. HTS-AT hierarchical architecture. The model consists of four stages of increasing embedding dimension: Stage 0 (2 blocks, 4 heads), Stage 1 (2 blocks, 8 heads), Stage 2 (6 blocks, 16 heads), and Stage 3 (2 blocks, 32 heads). PatchMerging is applied after Stages 0, 1, and 2, halving the spatial side length and doubling the feature dimension at each transition. Stage 3 has no PatchMerging. The input token sequence has spatial size $64 \times 64 = 4,096$ tokens with embedding dimension $D_0 = 96$; after three PatchMerging operations the final spatial size is $8 \times 8 = 64$ tokens with $D_3 = 768$.

Table 4. HTS-AT stage-level architectural parameters. The per-head dimension $d_h = D_\ell / H_\ell = 24$ is constant across all stages. S_ℓ is the spatial side length of the token grid at stage ℓ ; $N_w^\ell = S_\ell^2 / M$ is the number of attention windows, with $M = w^2 = 64$ tokens per window ($w = 8$).

Stage ℓ	Blocks B_ℓ	Heads H_ℓ	Dim D_ℓ	Spatial S_ℓ^2	Windows N_w^ℓ
0	2	4	96	64×64	64
1	2	8	192	32×32	16
2	6	16	384	16×16	4
3	2	32	768	8×8	1
Total heads H_{tot}		$2 \cdot 4 + 2 \cdot 8 + 6 \cdot 16 + 2 \cdot 32 = 8 + 16 + 96 + 64 = 184$			

Even-indexed blocks use standard Window Multi-head Self-Attention (W-MSA); odd-indexed blocks use Shifted-Window MSA (SW-MSA), which shifts the partition by $(w/2, w/2)$ tokens to enable cross-window interactions. At Stage 3, where $S_3 = w = 8$ so the grid coincides with a single window, the shift degenerates to zero and all blocks use W-MSA. The MLP sub-layer is a two-layer feed-forward network with hidden dimension $4D_\ell$ and GELU activation.

WindowAttention in detail. The W-MSA module at block (ℓ, b) first partitions the S_ℓ^2 tokens into $N_w^\ell = S_\ell^2 / M$ non-overlapping windows of $M = 64$ tokens each, then applies multi-head attention independently within each window. Queries, keys, and values are computed via a single fused projection $W_{\ell,b}^{QKV} \in \mathbb{R}^{D_\ell \times 3D_\ell}$:

$$\text{LN}(\mathbf{Z}^{(\ell,b-1)}) W_{\ell,b}^{QKV} \xrightarrow{\text{split}(D_\ell)} \mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N_w^\ell \times M \times D_\ell}. \quad (18)$$

Each tensor is then reshaped to isolate the H_ℓ heads, yielding per-head slices $\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h \in \mathbb{R}^{N_w^\ell \times M \times d_h}$ with $d_h = D_\ell / H_\ell = 24$.

For each head $h \in \{1, \dots, H_\ell\}$, attention is computed as:

$$\mathbf{H}_{\ell,b,h} = \text{Softmax} \left(\frac{\mathbf{Q}_h \mathbf{K}_h^\top}{\sqrt{d_h}} + \mathbf{B}_{\ell,b,h} \right) \mathbf{V}_h \in \mathbb{R}^{N_w^\ell \times M \times d_h}, \quad (19)$$

where $\mathbf{B}_{\ell,b,h} \in \mathbb{R}^{M \times M}$ is the learned relative position bias for head h of block (ℓ, b) . This bias encodes the relative spatial offset between each pair of tokens *within* a window; it is shared across all N_w^ℓ windows (broadcast) and is independent for each head. Concretely, it is read from a parameter table of shape $((2w-1)^2, H_\ell) = (225, H_\ell)$ via a fixed index mapping, so each head h has its own column of learnable values.

All H_ℓ head outputs are then concatenated and passed through the block-specific output projection $W_{\ell,b}^O \in \mathbb{R}^{D_\ell \times D_\ell}$ with bias $\mathbf{b}_{\ell,b}^O \in \mathbb{R}^{D_\ell}$, giving the W-MSA output

$$\mathbf{A}_{\ell,b} = [\mathbf{H}_{\ell,b,1} \parallel \dots \parallel \mathbf{H}_{\ell,b,H_\ell}] W_{\ell,b}^O + \mathbf{b}_{\ell,b}^O \in \mathbb{R}^{N_w^\ell \cdot M \times D_\ell}, \quad (20)$$

which is added to the residual stream in Eq. (16). For a complete step-by-step account of all tensor shapes across every stage, see <https://github.com/lorenzo-arcioni/ResiDual-CLAP/blob/main/README.md>.

A.3. Per-Head Residual Decomposition

The additive residual structure of Eq. (16) allows us to decompose the W-MSA contribution $\mathbf{A}_{\ell,b}$ exactly into independent per-head terms. Denoting by $W_{\ell,b,h}^O \in \mathbb{R}^{d_h \times D_\ell}$ the row slice of $W_{\ell,b}^O$ corresponding to head h (rows $[(h-1)d_h, hd_h)$), we distribute the output projection over heads:

$$\begin{aligned} \mathbf{A}_{\ell,b} &= [\mathbf{H}_{\ell,b,1} \parallel \cdots \parallel \mathbf{H}_{\ell,b,H_\ell}] W_{\ell,b}^O + \mathbf{b}_{\ell,b}^O \\ &= \sum_{h=1}^{H_\ell} \left(\mathbf{H}_{\ell,b,h} W_{\ell,b,h}^O + \frac{\mathbf{b}_{\ell,b}^O}{H_\ell} \right) = \sum_{h=1}^{H_\ell} \hat{\mathbf{H}}_{\ell,b,h}, \quad (21) \end{aligned}$$

where we define the *per-head projected contribution*

$$\hat{\mathbf{H}}_{\ell,b,h} = \mathbf{H}_{\ell,b,h} W_{\ell,b,h}^O + \frac{\mathbf{b}_{\ell,b}^O}{H_\ell} \in \mathbb{R}^{N_w^\ell \times M \times D_\ell}. \quad (22)$$

This decomposition is exact and follows from the linearity of matrix multiplication: because $W_{\ell,b}^O$ acts on the concatenation of all head outputs, each head h effectively multiplies only its own row slice $W_{\ell,b,h}^O$, and the output bias can be distributed equally among heads without any approximation.

The decomposition holds within a single stage, where D_ℓ is constant. Across stage boundaries, PatchMerging changes both spatial resolution and embedding dimension, breaking any global additive structure; cross-stage comparisons must therefore account for this.

Spatial aggregation. For each audio sample, stage, block, and head, we define a spatially aggregated scalar representation by averaging $\hat{\mathbf{H}}_{\ell,b,h}$ over all windows and token positions:

$$\hat{\mathbf{r}}_{\ell,b,h} = \frac{1}{N_w^\ell M} \sum_{i=1}^{N_w^\ell} \sum_{j=1}^M \hat{\mathbf{H}}_{\ell,b,h}[i, j, :] \in \mathbb{R}^{D_\ell}. \quad (23)$$

For the pre-projection analysis (Appendix B), the analogous aggregation is applied to the raw head output $\mathbf{H}_{\ell,b,h}$ before $W_{\ell,b}^O$:

$$\mathbf{r}_{\ell,b,h} = \frac{1}{N_w^\ell M} \sum_{i=1}^{N_w^\ell} \sum_{j=1}^M \mathbf{H}_{\ell,b,h}[i, j, :] \in \mathbb{R}^{d_h}. \quad (24)$$

Stacking these vectors across the n audio samples in the dataset yields the matrices $\hat{\mathbf{R}}_{\ell,b,h} \in \mathbb{R}^{n \times D_\ell}$ and $\mathbf{R}_{\ell,b,h} \in \mathbb{R}^{n \times d_h}$, which are the primary objects of our analysis. Because D_ℓ varies across stages, the ambient dimension of $\hat{\mathbf{R}}_{\ell,b,h}$ differs between stages, and any cross-stage comparison of dimensionality metrics must account for this varying ceiling.

B. Pre-Projection Head Analysis

The main analysis operates on *projected* head contributions $\hat{\mathbf{H}}_{\ell,b,h}$, which reflect each head’s influence on the residual stream. Here we describe a complementary analysis of *raw pre-projection* outputs $\mathbf{H}_{\ell,b,h}$, which characterise the intrinsic computational geometry of each head independently of W^O .

B.1. Motivation and Methodological Differences

The raw head output $\mathbf{H}_{\ell,b,h} \in \mathbb{R}^{N_w \times M \times d_h}$ is produced by the attention mechanism—specifically the weighted sum of value vectors—before any mixing across heads. It lives in a constant $d_h = 24$ -dimensional native space, offering two analytical advantages:

1. **Cross-stage comparability without projection.** $d_h = 24$ is identical for all 184 heads regardless of stage, so dimensionality metrics are directly comparable without mapping to an external space. This isolates head-internal geometry from the CLAP projection.
2. **Absence of W^O distortion.** W^O is a linear map that mixes head contributions and can rotate or rescale the geometry. The pre-projection space reflects what the attention mechanism *computes*; the post-projection space reflects what it *communicates* to the residual stream.

B.2. Extraction Procedure

The same forward hooks are reused. Instead of applying W_h^O , we directly aggregate $\mathbf{H}_{\ell,b,h}$ in the native d_h -dimensional space:

$$\tilde{\mathbf{r}}_{\ell,b,h} = \frac{1}{N_w M} \sum_{i=1}^{N_w} \sum_{j=1}^M \mathbf{H}_{\ell,b,h}[i, j, :] \in \mathbb{R}^{24}. \quad (25)$$

For each head we collect $\tilde{\mathbf{R}}_{\ell,b,h} \in \mathbb{R}^{n \times 24}$ and apply the same estimators of Section 3.3, with 24 as the maximum possible PCA dimension. No W^O or CLAP projection is applied.

B.3. Dimensionality Estimators in the Pre-Projection Space

All estimators from Section 3.3 apply with $\tilde{\mathbf{R}}_{\ell,b,h}$ replacing $\mathbf{R}_{\ell,b,h}$ and $d_h = 24$ as the ambient dimension ceiling. Key consequences:

PCA. The covariance $\tilde{\mathbf{C}} \in \mathbb{R}^{24 \times 24}$ has at most 24 non-zero eigenvalues. $d_{\text{PCA}}(\alpha) \leq 24$ for all stages, making the metric a direct measure of what fraction of the 24-dimensional native capacity each head exploits.

Table 5. Summary of notation used throughout the paper.

Symbol	Definition	Values / Notes
$\ell \in \{0, 1, 2, 3\}$	Stage index	
$b \in \{1, \dots, B_\ell\}$	Block index within stage ℓ	$B_\ell \in \{2, 2, 6, 2\}$
$h \in \{1, \dots, H_\ell\}$	Attention head index within block (ℓ, b)	$H_\ell = 4 \cdot 2^\ell \in \{4, 8, 16, 32\}$
S_ℓ	Spatial side length of the token grid at stage ℓ	$S_\ell \in \{64, 32, 16, 8\}$
$w = 8$	Attention window side length	
$M = w^2 = 64$	Tokens per attention window	
$N_w^\ell = S_\ell^2 / M$	Number of attention windows at stage ℓ	$N_w^\ell \in \{64, 16, 4, 1\}$
$d_h = 24$	Per-head feature dimension (constant across stages)	$d_h = D_\ell / H_\ell$
$D_\ell = H_\ell \cdot d_h$	Total embedding dimension at stage ℓ	$D_\ell \in \{96, 192, 384, 768\}$
$d_{\text{proj}} = 1024$	CLAP joint embedding dimension	
$W_{\ell,b}^{QKV} \in \mathbb{R}^{D_\ell \times 3D_\ell}$	Fused QKV projection at block (ℓ, b)	
$W_{\ell,b}^O \in \mathbb{R}^{D_\ell \times D_\ell}$	Output projection at block (ℓ, b)	
$W_{\ell,b,h}^O \in \mathbb{R}^{d_h \times D_\ell}$	Row slice of $W_{\ell,b}^O$ for head h	rows $[(h-1)d_h, hd_h)$
$\mathbf{b}_{\ell,b}^O \in \mathbb{R}^{D_\ell}$	Output projection bias at block (ℓ, b)	
$\mathbf{B}_{\ell,b,h} \in \mathbb{R}^{M \times M}$	Relative position bias for head h at block (ℓ, b)	
$\mathbf{H}_{\ell,b,h} \in \mathbb{R}^{N_w^\ell \times M \times d_h}$	Raw head output at block (ℓ, b) , head h (pre- $W_{\ell,b}^O$)	
$\hat{\mathbf{H}}_{\ell,b,h} \in \mathbb{R}^{N_w^\ell \times M \times D_\ell}$	Per-head projected contribution (post- $W_{\ell,b}^O$); see Eq. (22)	
$\mathbf{r}_{\ell,b,h} \in \mathbb{R}^{d_h}$	Spatially aggregated raw head output; see Eq. (24)	
$\hat{\mathbf{r}}_{\ell,b,h} \in \mathbb{R}^{D_\ell}$	Spatially aggregated projected head output; see Eq. (23)	
$\mathbf{R}_{\ell,b,h} \in \mathbb{R}^{n \times d_h}$	Dataset matrix of $\mathbf{r}_{\ell,b,h}$ across n samples	
$\hat{\mathbf{R}}_{\ell,b,h} \in \mathbb{R}^{n \times D_\ell}$	Dataset matrix of $\hat{\mathbf{r}}_{\ell,b,h}$ across n samples	
$P : \mathbb{R}^{768} \rightarrow \mathbb{R}^{1024}$	CLAP audio projection head (two-layer MLP with residual)	
n	Number of audio samples in the dataset	

Table 6. Comparison of the two analysis pipelines.

	Main (post- W^O)	Appendix (pre- W^O)
Representation	$\hat{\mathbf{H}}_{\ell,b,h}$	$\mathbf{H}_{\ell,b,h}$
Ambient dim	$D_\ell \in \{96, 192, 384, 768\}$	$d_h = 24$ (constant)
Analysis space	\mathbb{R}^{1024} (after P)	\mathbb{R}^{24} (native)
Hook location	Before .reshape, then $\times W_h^O$	Before .reshape only
W^O applied	Yes (via W_h^O)	No
P applied	Yes	No
Captures	Contribution to residual stream	Internal attention computation

L/N Ratio. With a fixed ambient dimension of 24, the ratio $d_{\text{PCA}_{99}}/d_{\text{TWO}_{\text{NN}}}$ isolates genuine manifold curvature from any ambient dimension effect, providing a cleaner nonlinearity diagnostic than in the post-projection case.

B.4. Interpretation and Relationship to Main Results

A head with low intrinsic dimensionality in the pre-projection space computes a low-rank attention pattern: the weighted combinations of value vectors collapse onto a small subspace of \mathbb{R}^{24} . Comparing pre- and post-projection results reveals the role of W^O : if dimensionality increases substantially after projection, W^O expands the represen-

tational geometry of that head in the residual stream; if it decreases, W^O compresses or mixes it.

Concretely, pre-projection analysis is best suited to studying individual head specialisation in isolation; post-projection analysis—the perspective adopted in the main body—is best suited to studying how heads collectively shape the residual stream and, ultimately, the CLAP embedding used for zero-shot classification.

C. Extended Dimensionality Analysis

This appendix provides comprehensive quantitative details and additional visualizations complementing the main results in Section 4.1.

C.1. Detailed Block-Level Statistics

Table 7 reports complete block-wise metrics across all 12 transformer blocks in HTS-AT, aggregating over the heads within each block as described in Section 3.3.

Relationship to Architecture. As illustrated in Figure 6, the spatial resolution decreases progressively through the network due to patch merging between stages. While this affects the number of tokens N processed by each attention

head, our analysis focuses on the intrinsic dimensionality of the *head dimension* $d_h = 24$ after spatial aggregation (Eq. 24). Thus, the reported metrics characterize the semantic complexity of head representations independent of spatial resolution effects.

The hierarchical structure creates natural breakpoints for dimensionality analysis:

- **Stage 1 (Blocks 0–1):** High spatial resolution ($T/2P \times F/2P$) but limited capacity ($D_0 = 96$). Early fusion of local spectral-temporal patterns.
- **Stage 2 (Blocks 2–3):** First dimensionality jump coincides with $2\times$ patch merging and head doubling. Transition from local to intermediate-scale features.
- **Stage 3 (Blocks 4–9):** Deepest stage with 6 blocks enables iterative refinement at fixed spatial scale ($T/4P \times F/4P$) and capacity ($D_2 = 384$). Gradual dimensionality growth reflects progressive feature abstraction.
- **Stage 4 (Blocks 10–11):** Maximum capacity ($D_3 = 768$) without further spatial reduction. Minimal dimensionality increase suggests saturation.

Table 7. Block-wise aggregated dimensionality metrics for TinySOL dataset. Blocks 0–1 (Stage 1), 2–3 (Stage 2), 4–9 (Stage 3), 10–11 (Stage 4). L = Linear ID ($d_{PCA_{99}}$), N = Nonlinear ID (TwoNN), L/N = Linear-nonlinear ratio, EVR1 = First PC variance explained.

Block	Stage	Heads	L	N	L/N	EVR1
0	1	4	3.75	3.94	0.95	0.878
1	1	4	8.00	4.94	1.62	0.575
2	2	8	12.75	6.15	2.07	0.403
3	2	8	17.50	6.75	2.59	0.460
4	3	16	18.00	6.93	2.60	0.388
5	3	16	20.13	7.08	2.84	0.279
6	3	16	21.25	7.40	2.87	0.250
7	3	16	21.38	7.42	2.88	0.243
8	3	16	22.25	8.37	2.66	0.217
9	3	16	21.88	8.11	2.70	0.241
10	4	32	22.25	8.49	2.62	0.262
11	4	32	21.59	8.00	2.70	0.272

Interpretation.

- Block 0 operates near the linear regime ($L/N \approx 1$), with almost 88% variance in the first PC, indicating extremely constrained early processing.
- The largest single-block jump occurs at the Stage 1→2 transition (blocks 1→2: $\Delta L = +4.75$, +59%), corresponding to doubling of attention heads (4→8) and hidden dimension (96→192).

- Stage 3 exhibits gradual linear ID growth (18.00 → 22.25 over 6 blocks) despite constant architecture, suggesting intra-stage feature refinement through depth.
- Stage 4 shows minimal progression (blocks 10→11: $\Delta L = -0.66$), consistent with representational saturation observed in the main text.

C.2. Extended Cross-Dataset Analysis

We replicate the full layer-wise analysis on ESC-50 and VocalSound to validate architectural generalizability. Tables 8 and 9 present complete statistics.

Table 8. Layer-wise dimensionality metrics for ESC-50 dataset (50 environmental sound classes, 1000 stratified samples).

Layer	$d_{PCA_{99}}$	TwoNN	PR	EVR1
L0	5.2 ± 2.1	4.2 ± 0.7	1.8 ± 0.8	0.741 ± 0.187
L1	14.3 ± 2.8	6.2 ± 0.6	3.9 ± 1.3	0.449 ± 0.109
L2	20.1 ± 2.0	7.4 ± 0.8	7.5 ± 1.8	0.281 ± 0.081
L3	20.3 ± 1.2	7.8 ± 1.1	7.4 ± 1.9	0.279 ± 0.074

Table 9. Layer-wise dimensionality metrics for VocalSound dataset (6 vocal sound categories, 1000 stratified samples).

Layer	$d_{PCA_{99}}$	TwoNN	PR	EVR1
L0	6.1 ± 2.3	4.6 ± 0.8	2.1 ± 1.0	0.698 ± 0.192
L1	16.2 ± 2.4	6.7 ± 0.5	4.5 ± 1.5	0.421 ± 0.117
L2	21.9 ± 1.9	8.0 ± 0.9	8.2 ± 2.0	0.263 ± 0.079
L3	22.7 ± 1.1	8.6 ± 0.8	8.3 ± 1.6	0.254 ± 0.071

Cross-Dataset Consistency Analysis. Despite differing semantic granularities (ESC-50: 50 classes, VocalSound: 6 classes, TinySOL: 14 classes), layer-wise trends remain remarkably stable:

- **L0 Concentration:** All datasets exhibit EVR1 $\geq 69\%$ in Stage 1, confirming universal early spectral concentration.
- **L1 Expansion:** The L0→L1 dimensionality jump is consistent (TinySOL: +9.2, ESC-50: +9.1, VocalSound: +10.1 for $d_{PCA_{99}}$), with coefficient of variation across datasets $CV = 0.06$.
- **L2-L3 Saturation:** All datasets show similar modest L2→L3 growth ($\Delta d < 2.0$), despite L3 having $2\times$ the heads of L2, indicating architecture-driven capacity limits.
- **L/N Ratio Convergence:** By Stage 4, all datasets reach $L/N \approx 2.6$ – 2.7 , suggesting a universal nonlinear complexity regime independent of semantic domain.

C.3. Statistical Validation

C.3.1. ANOVA RESULTS

One-way ANOVA tests for layer differences on TinySOL dataset:

Table 10. Statistical significance of layer effects on dimensionality metrics (TinySOL, $n = 184$ heads). All tests use $\alpha = 0.05$.

Metric	F-statistic	p-value	Significance
$d_{\text{PCA}_{99}}$	262.64	< 0.001	***
TwoNN	58.93	< 0.001	***
PR	44.74	< 0.001	***
EffRank	76.03	< 0.001	***
EVR(PC1)	118.47	< 0.001	***

C.3.2. POST-HOC PAIRWISE COMPARISONS

Bonferroni-corrected pairwise t-tests for $d_{\text{PCA}_{99}}$ (6 comparisons, $\alpha_{\text{corrected}} = 0.0083$):

Table 11. Pairwise layer comparisons for linear intrinsic dimensionality (TinySOL). All comparisons significant at corrected $\alpha = 0.0083$.

Comparison	$\Delta d_{\text{PCA}_{99}}$	Cohen’s d	p-value
L0 vs L1	9.22	3.86	< 0.001
L0 vs L2	14.93	7.12	< 0.001
L0 vs L3	16.04	8.94	< 0.001
L1 vs L2	5.71	2.34	< 0.001
L1 vs L3	6.82	3.19	< 0.001
L2 vs L3	1.11	0.78	< 0.001

All effect sizes exceed Cohen’s threshold for ”large” effects ($d > 0.8$), with the L0 vs L3 comparison exhibiting extremely large effects ($d > 8$), confirming substantial representational differences across layers.

C.4. Additional Visualizations

C.4.1. PC1 DOMINANCE AND BOXPLOTS

Figure 7 presents complementary views of dimensionality structure.

Observations from Panel (a):

- Only 2 heads in L0 (2.1%) fall below the 50% EVR1 threshold, compared to 89% of L3 heads, demonstrating near-universal early concentration.
- The EVR1 distribution shifts from unimodal (L0: concentrated near 0.7–0.8) to bimodal (L3: peaks at 0.2–0.3 and 0.35–0.4), suggesting emergence of head sub-populations with distinct specialization levels.

Observations from Panel (b):

- PR and EffRank exhibit parallel scaling, confirming their measurement of related spectral properties.
- TwoNN shows compressed scale relative to PCA99, visually emphasizing the linear-nonlinear gap discussed in the main text.
- Outliers (marked as individual points beyond whiskers) are rare in L0-L1 but frequent in L2-L3, consistent with increased head-level heterogeneity.

C.5. Computational Details

All analyses were performed on an NVIDIA A100 GPU (40GB) using PyTorch 2.0.1 and Python 3.10. Key implementation details:

- **Head Extraction:** Forward hooks registered via `torch.nn.Module.register_forward_hook`. Batch size 100 for extraction to balance memory and throughput.
- **PCA:** Computed via `sklearn.decomposition.PCA` with full SVD solver. Eigenvalue thresholding at machine epsilon ($\sim 10^{-7}$) to remove numerical noise.
- **TwoNN:** Implemented using `skdim.id.TwoNN` with default parameters (no k selection required).
- **MLE:** `skdim.id.MLE` with $k = 20$ neighbors, standard Euclidean metric.
- **Statistical Tests:** `scipy.stats` functions (`f_oneway`, `ttest_ind`, `spearmanr`) with standard settings.

Total extraction time: ~ 45 minutes per dataset (1000 samples \times 184 heads). Analysis pipeline code available at [https://github.com/\[ANONYMOUS\]/residual-audio](https://github.com/[ANONYMOUS]/residual-audio).

C.6. Reproducibility Checklist

To facilitate replication:

- Random seeds: 42 (Python), 42 (NumPy), 42 (PyTorch)
- CLAP version: `laion/clap-htsat-unfused` checkpoint from HuggingFace
- Audio preprocessing: CLAP default (64-band mel, 10s duration, 48kHz resampling)



Figure 7. Extended dimensionality analysis. (a) First principal component variance explained across all 184 heads. Horizontal line at 50% marks equal-contribution threshold. Sharp decline from L0 (mean 73%) to L3 (mean 27%) quantifies transition from low-rank to distributed representations. (b) Boxplot comparison of four key metrics across layers, revealing consistent monotonic trends and increasing intra-layer variance in deeper stages (note wider boxes for L2-L3).

- Dataset versions: ESC-50 v2.0, TinySOL v3.0, Vocal-Sound official release
- Stratified sampling:
`sklearn.model_selection.StratifiedShuffleSplit`
with 1000 samples