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

February 19, 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 main components: (1) comprehensive analysis of the residual stream structure in CLAP’s audio encoder to identify specialization patterns, and (2) implementation of spectral reweighting techniques to enhance task performance. This section describes the analysis methodology, with implementation details of the ResiDual adaptation deferred to Section 3.5.

3.1. Model Architecture and Notation

We analyze the HTS-AT (Hierarchical Token-Semantic Audio Transformer) architecture (?), which serves as the audio encoder in CLAP (?). HTS-AT processes audio through four hierarchical stages with depths [2, 2, 6, 2], employing Swin-Transformer blocks with window-based self-attention of size $w = 8$.

Figure 1 illustrates the complete HTS-AT pipeline. Input spectrograms undergo progressive spatial downsampling via patch merging between stages, while the number of attention heads doubles at each transition, following a standard hierarchical vision transformer design adapted for audio.

The architecture parameters are:

- Embedding dimension: $d_{\text{emb}} = 96$
- Number of attention heads per stage: $H_\ell = 4 \cdot 2^\ell$, $\ell = 0, 1, 2, 3$ (i.e., $(H_0, H_1, H_2, H_3) = (4, 8, 16, 32)$)
- Head dimension: $d_h = 24$ (constant across layers)
- Total layer capacity: $D_\ell = H_\ell \cdot d_h$, yielding $D_0 = 96$, $D_1 = 192$, $D_2 = 384$, $D_3 = 768$
- Total attention heads: $H_{\text{tot}} = \sum_{\ell=0}^3 (H_\ell \cdot \text{depth}_\ell) = 184$

For each attention head h in layer ℓ and block b , we denote the head output *before* the output projection as

$$\mathbf{H}_{\ell,b,h} \in \mathbb{R}^{N \times M \times d_h}, \quad (1)$$

where N is the number of spatial windows and $M = w^2$ is the number of patches per window. This corresponds to the weighted value vectors:

$$\mathbf{H}_{\ell,b,h} = \text{Attn}_h(\mathbf{Q}, \mathbf{K})\mathbf{V}_h, \quad (2)$$

with $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{N \times M \times d_h}$ as the query and key projections, and $\mathbf{V}_h \in \mathbb{R}^{N \times M \times d_h}$ as the value projections specific to head h .

3.2. Residual Stream Extraction

To analyze the intrinsic structure of attention head representations, we extract pre-projection head outputs using forward hooks registered on each attention module. Specifically, for each head, we capture $\mathbf{H}_{\ell,b,h}$ before it undergoes the final linear transformation W^O .

Aggregation Strategy. Since audio spectrograms produce variable-length token sequences depending on window partitioning, we aggregate head outputs spatially by mean pooling over both the window and patch dimensions:

$$\mathbf{r}_{\ell,b,h} = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \mathbf{H}_{\ell,b,h}[i, j, :] \in \mathbb{R}^{d_h}. \quad (3)$$

This yields a single d_h -dimensional representation per audio sample per head, enabling cross-sample analysis while preserving head-specific characteristics. Each $\mathbf{r}_{\ell,b,h}$ can then be used for downstream analysis or aggregation across heads and blocks.

Note: This spatial aggregation is motivated by two considerations: (i) the computational cost of analyzing the full set of tokens would be prohibitive, and (ii) averaging over windows and patches provides a concise, holistic representation of each head’s behavior for a given audio sample.

Dataset Sampling. We extract representations from stratified samples across three audio classification benchmarks:

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

For our analysis, we collect all the samples from each dataset. Audio preprocessing follows CLAP’s standard pipeline: 64-band mel-spectrogram with 10-second padding/truncation.

3.3. Intrinsic Dimensionality Analysis

To characterize the effective complexity of head representations, we employ a multi-faceted dimensionality analysis combining linear and nonlinear estimators.

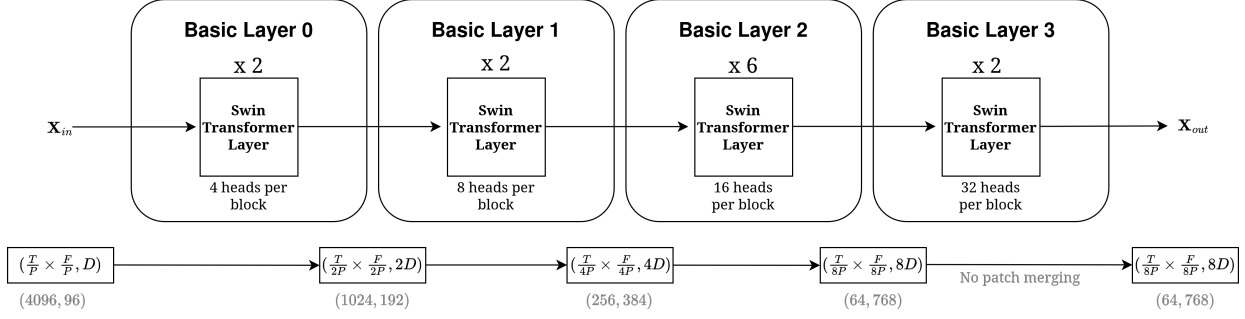


Figure 1. HTS-AT hierarchical architecture. The model consists of four basic layers (stages) with increasing complexity: Stage 1 (2 blocks \times 4 heads), Stage 2 (2 blocks \times 8 heads), Stage 3 (6 blocks \times 16 heads), and Stage 4 (2 blocks \times 32 heads). Patch merging between stages reduces spatial resolution while doubling the feature dimension. Input spectrogram dimensions are $T/P \times F/P = 64 \times 64 = 4096$ patches with $D = 96$ channels. The final output has spatial size $(T/8P) \times (F/8P) = 64 \times 768$ before global pooling. Note that Stage 4 omits patch merging to preserve spatial resolution for fine-grained modeling.

3.3.1. LINEAR DIMENSIONALITY ESTIMATORS

PCA-based Metrics. For each head, we compute the covariance matrix $\mathbf{C}_{\ell,b,h} = \frac{1}{n-1} \mathbf{R}_{\ell,b,h}^\top \mathbf{R}_{\ell,b,h}$, where $\mathbf{R}_{\ell,b,h} \in \mathbb{R}^{n \times d_h}$ stacks all aggregated representations. Eigendecomposition yields ordered eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{d_h}$.

We define linear intrinsic dimensionality as the minimum number of principal components capturing α variance:

$$d_{\text{PCA}}(\alpha) = \arg \min_k \left\{ \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^{d_h} \lambda_i} \geq \alpha \right\} \quad (4)$$

We evaluate $\alpha \in \{0.90, 0.95, 0.99\}$ to capture both coarse and fine-grained dimensionality.

Participation Ratio. The Participation Ratio (PR) (?) quantifies spectral dispersion:

$$\text{PR} = \frac{(\sum_{i=1}^{d_h} \lambda_i)^2}{\sum_{i=1}^{d_h} \lambda_i^2} \quad (5)$$

High PR indicates uniform variance distribution across dimensions, while low PR suggests concentration on few dominant directions.

Effective Rank. Based on the normalized eigenspectrum $p_i = \lambda_i / \sum_j \lambda_j$, we compute Shannon entropy:

$$H = - \sum_{i=1}^{d_h} p_i \log p_i \quad (6)$$

The Effective Rank is then:

$$\text{EffRank} = \exp(H) \quad (7)$$

This represents the equivalent dimension of a uniform distribution with identical entropy, providing a continuous measure of effective dimensionality.

3.3.2. NONLINEAR DIMENSIONALITY ESTIMATORS

To complement linear analyses, we also employ non-linear intrinsic dimensionality estimators that capture potential manifold curvature in the head representations.

TwoNN Estimator. The TwoNN method (?) estimates intrinsic dimensionality from the empirical distribution of nearest-neighbor distance ratios. For each pooled head vector $\mathbf{r}_{\ell,b,h} \in \mathbb{R}^{d_h}$, let $r_1^{(i)}$ and $r_2^{(i)}$ denote the Euclidean distances to its first and second nearest neighbors within the dataset of all samples for that head. The estimator is:

$$d_{\text{TwoNN}} = \left(\frac{1}{n} \sum_{i=1}^n \log \frac{r_2^{(i)}}{r_1^{(i)}} \right)^{-1}, \quad (8)$$

where n is the number of samples. This nonparametric approach is robust to manifold curvature and does not assume linear structure.

MLE Estimator. The Maximum Likelihood Estimator (MLE) (?) assumes local uniformity of the data-generating distribution. For each pooled vector $\mathbf{r}_{\ell,b,h}$, we consider its k nearest neighbors and compute:

$$\hat{d}_{\text{MLE}}(\mathbf{r}_{\ell,b,h}) = \left(\frac{1}{k-1} \sum_{j=1}^{k-1} \log \frac{r_k(\mathbf{r}_{\ell,b,h})}{r_j(\mathbf{r}_{\ell,b,h})} \right)^{-1}, \quad (9)$$

where $r_j(\mathbf{r}_{\ell,b,h})$ is the distance to the j -th nearest neighbor. The global MLE estimate for each head is obtained by averaging over all samples. We set $k = 20$, following standard practice.

Implementation Details. For each head, we first compute a single PCA decomposition of the pooled vectors $\mathbf{r}_{\ell,b,h}$ to obtain the covariance structure and the eigenvalues λ_i . From these, we derive the Participation Ratio, Effective Rank, and PCA-based intrinsic dimensionalities (e.g., number of components capturing 90%, 95%, 99% variance). The TwoNN and MLE estimators are then applied directly to the same set of pooled vectors.

All computations are performed for every head individually, and the resulting metrics are stored along with meta-data (layer, block, head index) to allow aggregation at the block or layer level for statistical analysis.

3.3.3. BLOCK-LEVEL AGGREGATION

HTS-AT organizes attention heads into *blocks* within each layer (Stage). To analyze coarser architectural patterns, we aggregate metrics at the block level. For block B containing heads \mathcal{H}_B , we compute:

$$\bar{m}_B = \frac{1}{|\mathcal{H}_B|} \sum_{h \in \mathcal{H}_B} m_h \quad (10)$$

for each metric $m \in \{\text{PCA}_{99}, \text{TwoNN}, \text{PR}, \text{EVR}_1\}$.

We define the **Linear-Nonlinear Ratio** as:

$$\text{Ratio}_B = \frac{\bar{d}_{\text{PCA}_{99}}}{\bar{d}_{\text{TwoNN}}} \quad (11)$$

Values near 1 suggest linear manifold structure, while higher values indicate nonlinear complexity beyond what PCA captures.

3.4. Statistical Analysis

To validate layer-wise progression and heterogeneity, we perform:

- **One-way ANOVA** to test for significant differences in dimensionality metrics across layers
- **Post-hoc pairwise comparisons** with Bonferroni correction ($\alpha = 0.05$)
- **Spearman rank correlation** to assess monotonic trends with layer depth
- **Effect size estimation** via Cohen’s d for layer-wise comparisons

All analyses use `scipy.stats` and `scikit-learn` implementations with random seed 42 for reproducibility.

3.5. ResiDual Spectral Reweighting

[TO BE COMPLETED: This section will describe:

- PCA decomposition of selected head outputs based on dimensionality analysis
- Spectral reweighting strategy (amplification of task-relevant components)
- Integration into HTS-AT forward pass
- Training/fine-tuning protocol for reweighted model
- Hyperparameter selection for component retention and scaling factors

]

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.

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

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

4.1.2. BLOCK-LEVEL ANALYSIS

Figure 2 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 1. Stage boundaries (blocks 1 \rightarrow 2, 3 \rightarrow 4, 9 \rightarrow 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 \rightarrow 2 doubles both heads (4 \rightarrow 8) and dimension (96 \rightarrow 192) yet achieves only modest dimensionality growth, while the Stage 2 \rightarrow 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 rep-

resentational manifold. Unlike earlier stages, capacity expansion does not translate into increased effective dimensionality, indicating a saturation of task-relevant feature complexity.

HTSAT Block-wise Metrics				
Block	L	N	L/N	EVR1
11	22.62	8.91	2.54	0.19
10	23.31	9.03	2.58	0.15
9	22.62	9.91	2.28	0.16
8	22.88	9.76	2.34	0.15
7	22.56	9.51	2.37	0.16
6	22.19	8.62	2.57	0.18
5	21.12	8.12	2.60	0.24
4	19.44	7.34	2.65	0.24
3	18.62	7.36	2.53	0.27
2	13.50	7.42	1.82	0.44
1	6.75	6.20	1.09	0.67
0	2.75	4.70	0.59	0.92

Figure 2. 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 ($d_{PCA_{99}}$), **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 lin-

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

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

4.1.4. INDIVIDUAL HEAD VARIABILITY

Figure 3 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 ($\text{EVR}(\text{PC1})$) 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 $\text{EVR}(\text{PC1})$ below 40% (Layer 2) and often under 20% (Layer 3), confirming that deeper representations spread information more evenly across

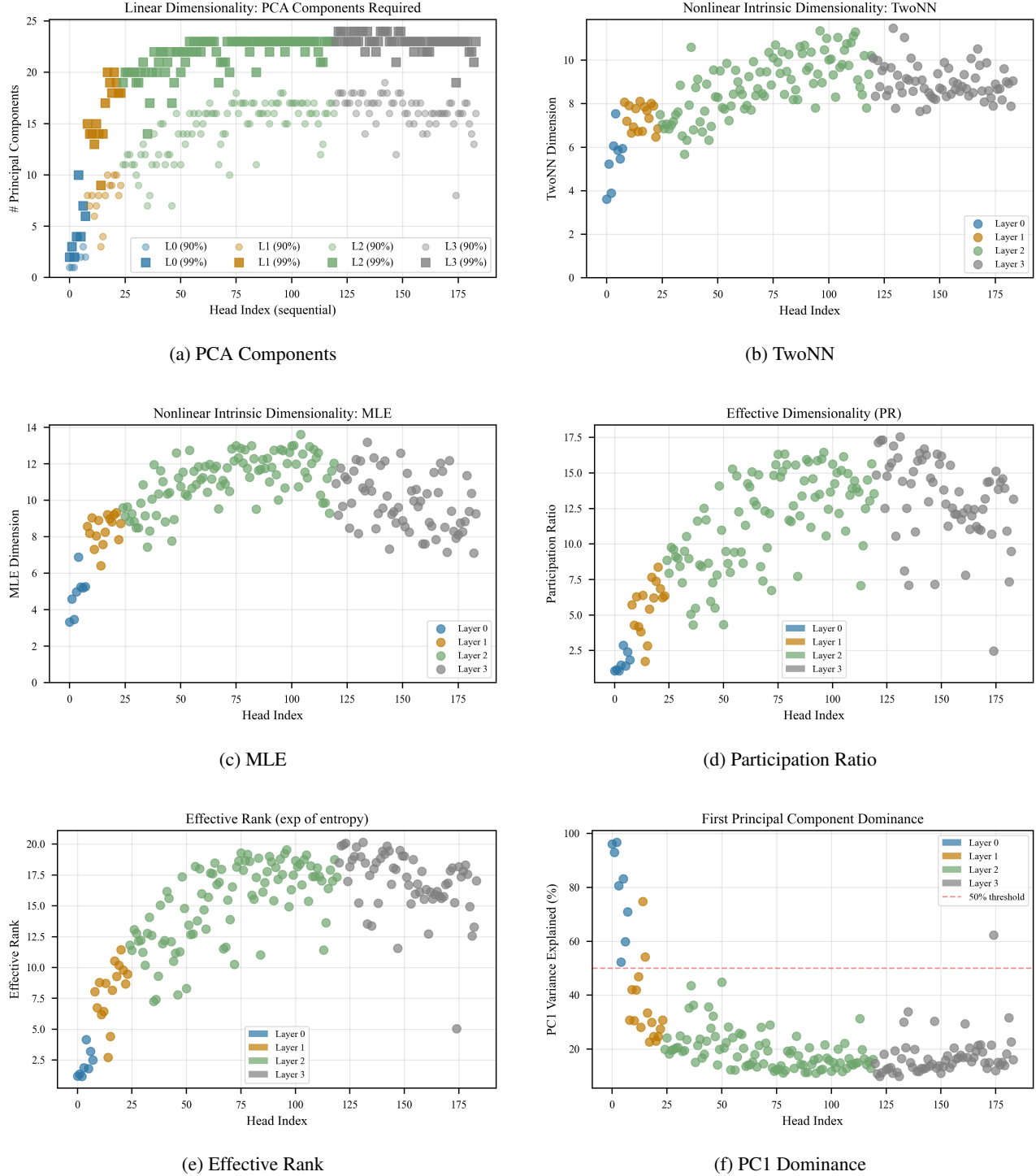


Figure 3. Head-level dimensionality analysis across HTS-AT layers.

multiple dimensions, consistent with increased representational richness and reduced linear redundancy.

Across layers, representational dimensionality increases and becomes more uniformly distributed across atten-

tion heads. Early layers exhibit highly constrained and heterogeneous representations, with PCA₉₉ 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 4 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.

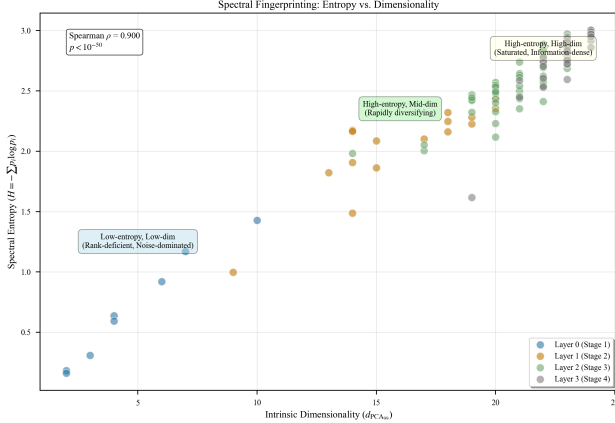


Figure 4. Spectral entropy vs. $d_{\text{PCA}_{99}}$ for all 184 HTS-AT heads on ESC-50 ($\rho = 0.900$, $p < 10^{-50}$). Color denotes layer.

4.2.2. CROSS-HEAD SIMILARITY STRUCTURE

Figure 5 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 bor-

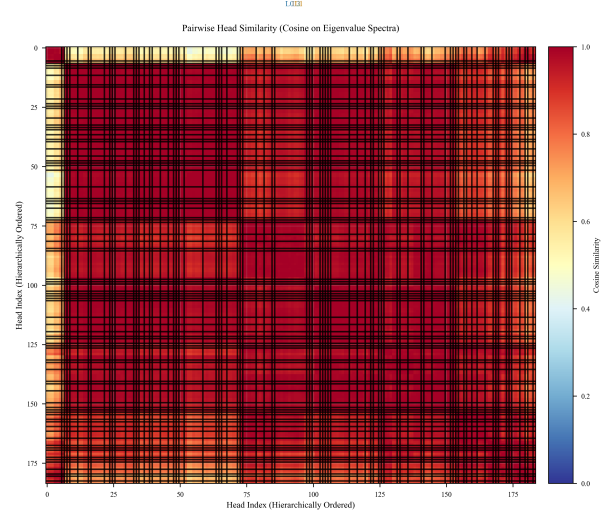


Figure 5. Pairwise head similarity (cosine on eigenvalue spectra, Ward ordering). The dominant low-similarity stripe corresponds to L0 heads, whose rank-deficient spectra are structurally distinct from all other layers.

der. This dichotomy indicates that the qualitative spectral transition occurs *at the Stage 1 \rightarrow 2 boundary*, after which all heads converge to broadly similar variance concentration profiles regardless of depth. The architectural capacity increases at later stage boundaries do not introduce analogous spectral discontinuities.

4.2.3. BLOCK-LEVEL ABLATION

Figure 6 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%).

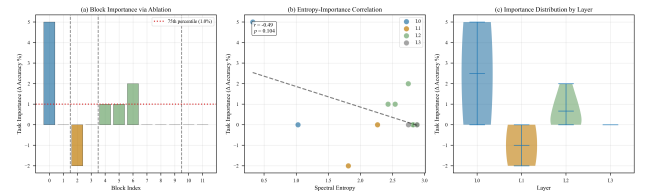


Figure 6. Block-level ablation on ESC-50 zero-shot classification. (a) Per-block Δ_b ; positive values indicate task-critical blocks. (b) Entropy vs. importance ($r = -0.49$, $p = 0.104$). (c) Layer-wise Δ_b distribution.

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:

1. **Layer Selection:** Given the sharp dimensionality increase at Stage 2 (Layer 1, $d_{\text{PCA}_{99}} = 16.1$) and continued expansion in Stage 3 (Layer 2, $d_{\text{PCA}_{99}} = 21.8$), 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.
2. **Component Retention:** For L2-L3 heads with $d_{\text{PCA}_{99}} \sim 22$, we can safely reduce to $k \approx 16$ components (73% of original) while retaining 99% variance, enabling efficient reweighting with minimal information loss.
3. **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.]**
4. **Head-Specific Strategies:** The intra-layer heterogeneity in middle layers (L1 CV = 0.19 for $d_{\text{PCA}_{99}}$) 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 ($d_{\text{PCA}_{99}}$: $5.9 \rightarrow 21.9$), 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→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 \hat{c} 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. ??) 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>.

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

- 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>.
- 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. Extended Dimensionality Analysis

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

A.1. Detailed Block-Level Statistics

Table 3 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 1, 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. ??). 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.

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.

Table 3. 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

- Stage 4 shows minimal progression (blocks 10→11: $\Delta L = -0.66$), consistent with representational saturation observed in the main text.

A.2. Extended Cross-Dataset Analysis

We replicate the full layer-wise analysis on ESC-50 and VocalSound to validate architectural generalizability. Tables 4 and 5 present complete statistics.

Table 4. 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 5. 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× 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.

A.3. Statistical Validation

A.3.1. ANOVA RESULTS

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

Table 6. Statistical significance of layer effects on dimensionality metrics (TinySOL, $n = 184$ heads). All tests use $\alpha = 0.05$.

Metric	F-statistic	p-value	Significance
$d_{PCA_{99}}$	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	***

A.3.2. POST-HOC PAIRWISE COMPARISONS

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

Table 7. Pairwise layer comparisons for linear intrinsic dimensionality (TinySOL). All comparisons significant at corrected $\alpha = 0.0083$.

Comparison	$\Delta d_{PCA_{99}}$	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.

A.4. Additional Visualizations

A.4.1. PC1 DOMINANCE AND BOXPLOTS

Figure 7 presents complementary views of dimensionality structure.



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).

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 subpopulations 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.

A.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).

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