

Bi-Orthogonal Factor Decomposition for Vision Transformers

Fenil R. Doshi^{*a,b} Thomas Fel^{*b} Talia Konkle^{a,b} George A. Alvarez^{a,b}

^aDept. of Psychology, Harvard University ^bKempner Institute, Harvard University

fenildoshi.com/bfd-transformers

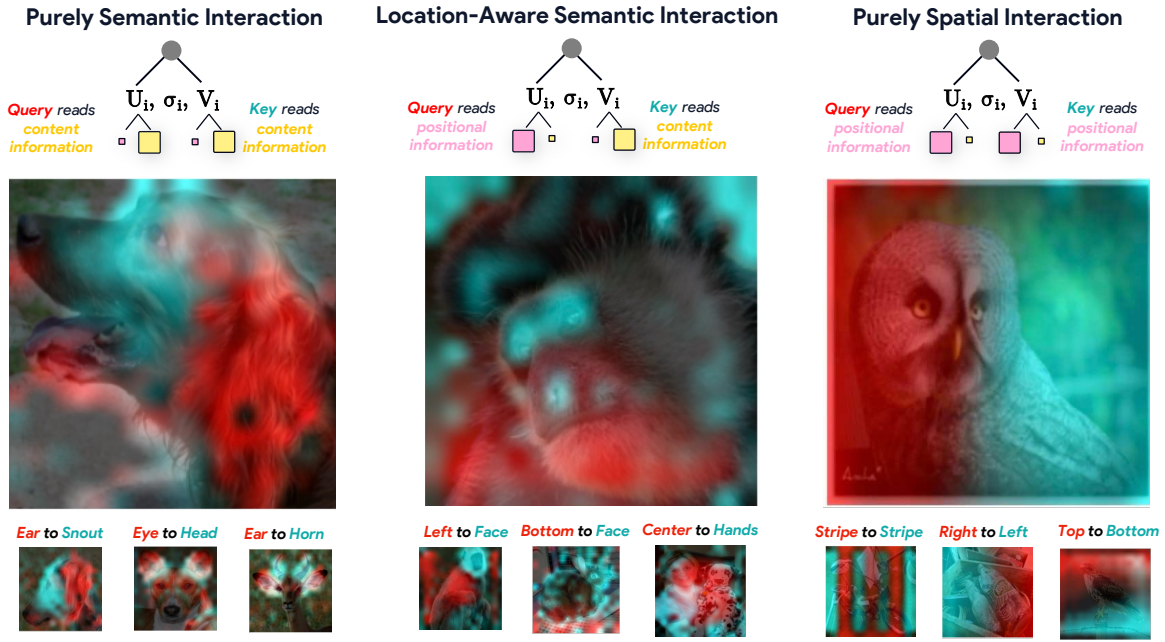


Figure 1. **Content and Position Interactions in DINOv2.** Representative bi-orthogonal modes in DINOv2 whose activation patterns illustrate three different interaction. Left: Content–content modes activate on semantic interactions such as object parts, revealing part-to-part or part-to-whole correspondences. Middle: Content–position modes show localization-aware semantic interactions, where semantic features are modulated by spatial context. Right: Position–position modes activate on purely spatial interaction, like left–right flows, top–bottom variations, and Fourier-like patterns, highlighting geometric organization without semantic selectivity. These examples show how different interaction types contribute to information flow in attention.

Abstract

Self-attention is the central computational primitive of Vision Transformers, yet we lack a principled understanding of what information attention mechanisms exchange between tokens. Attention maps describe where weight mass concentrates; they do not reveal whether queries and keys trade position, content, or both. We introduce Bi-orthogonal Factor Decomposition (BFD), a two-stage analytical framework: first, an ANOVA-based decomposition statistically disentangles token activations into orthogonal positional

and content factors; second, SVD of the query-key interaction matrix QK^T exposes bi-orthogonal modes that reveal how these factors mediate communication. After validating proper isolation of position and content, we apply BFD to state-of-the-art vision models and uncover three phenomena. (i) Attention operates primarily through content. Content-content interactions dominate attention energy, followed by content-position coupling. DINOv2 allocates more energy to content-position than supervised models and distributes computation across a richer mode spectrum. (ii) Attention mechanisms exhibit specialization: heads differentiate into content-content, content-position, and position-position operators, while singular modes within heads show analogous

* Denotes equal contribution.

specialization. (iii) DINOv2’s superior holistic shape processing emerges from intermediate layers that simultaneously preserve positional structure while contextually enriching semantic content.

Overall, BFD exposes how tokens interact through attention and which informational factors – positional or semantic – mediate their communication, yielding practical insights into vision transformer mechanisms.

1. Introduction

Vision Transformers (ViTs) [22, 27, 62, 81] have become the workhorses of modern computer vision, powering recognition systems deployed in medical and other safety-critical settings [4, 83] and serving as the backbone of large-scale generative models [92, 94] used by millions [16]. Understanding these systems is therefore both a practical imperative and a scientific opportunity [51, 70]. From a scientific perspective, ViTs are among the most scalable [3, 23, 90] and capable vision models available [13, 41, 77, 91]; it is natural to ask whether their internal computations share structure with human visual processing [17, 31, 47, 50, 57, 76] and, if so, where such similarities emerge [12, 19, 46, 54]? To answer those question, we turn to the interpretability field [26, 39, 42], which has proposed a range of theories [21, 29, 65] and methods [10, 44, 89] for understanding model internals. Early interpretability efforts centered on attribution [30, 36, 58, 66, 69, 74, 89]. These methods primarily aim to reveal where a model attends or which regions of an input most influence its prediction. While such tools offer valuable intuition, they face well-documented limitations [1, 18, 45, 59, 72, 73]. They often highlight broad image regions without clarifying what visual patterns drive the decision, nor how these signals are combined internally. In response, concept-based methods emerged [28, 33, 34, 38, 44, 48, 64, 71, 82, 93], shifting focus from spatial attribution to semantic characterization. Rather than showing *where* the model looks, these approaches aim to identify *what* it has learned – isolating human-interpretable “visual atoms” [70, 79] or latent concepts within its representations [10]. Concept activation vectors [44], network dissection [8] or feature visualization [32, 40, 60] exemplify this paradigm. Such analyses have revealed, for instance, that deep networks spontaneously form units selective for texture, object parts, or higher-order scene semantics [35, 49]. Yet, despite their insights, these techniques remain representation-centric: they explain *what* information is encoded, not *how* that information arises. They stop short of elucidating the circuits [53] or computations [67] that produce these representations – a question that has become central to modern interpretability.

To uncover these internal computations, attention mechanisms present a natural entry point. Self-attention is a core

component of ViTs as it mediates communication between tokens and thus provides an explicit interface for information exchange. Existing analyses have mostly interpreted attention through its spatial patterns – visualizing which tokens attend to which others across layers [14, 87]. However, these maps reflect only the surface structure of interactions: they indicate the amount of attention exchanged, not the content of that exchange or the reasons for it [2, 7, 84]. This gap has fueled growing skepticism over whether raw attention maps genuinely constitute explanations [9, 78]. Recently, however, a more mechanistic turn has emerged. Rather than treating attention weights as explanatory endpoints, researchers probe the underlying query-key interactions directly via spectral structure [63]. These methods aim to identify dominant interaction modes and relate them to concrete behaviors, echoing circuit-level analyses pioneered in large language models [11, 15, 75, 80] (e.g., induction heads [61], copying mechanisms [55], compositional reasoning [5, 52]). Yet, even when dominant modes are identified, their nature often remains ambiguous: are these modes driven primarily by positional geometry, by semantic content, or by global trends such as the mean activation direction? Without a principled factor attribution framework [75] (one that disentangles these potential sources) it remains unclear whether the apparent contextual integration observed in deeper layers reflects genuinely content-driven computation or merely the imprint of positional structure. This motivates our first question.

Question 1. *What informational factors – positional, semantic, or global – mediate token communication through attention, and how does this composition vary across layers and between supervised vs. self-supervised models?*

Understanding what information flows between tokens naturally leads to asking *how* this flow is organized at the network level. The architecture of multi-head attention introduces a second source of ambiguity here as each head could, in principle, implement a distinct computational role (tracking spatial relations, aggregating semantics, ...) supporting a modular view of attention. Alternatively, information may be distributed across heads, with individual heads contributing only partial, entangled signals to a collective computation. These contrasting modular versus distributed hypotheses carry deep implications for interpretability: the former supports circuit-level mechanistic explanations at the head level, while the latter implies that per-head analyses may be fundamentally incomplete. This leads to our second question.

Question 2. *Do attention heads and their constituent singular modes exhibit functional specialization into distinct informational operators (content-content, content-position, position-position), and is this specialization consistent across layers and models?*

Understanding the organization of attention naturally raises a further question: what determines it? If head specializations exist, are they intrinsic to the architecture or induced by the learning paradigm? Evidence suggests that training objectives and data statistics strongly influence representational geometry and inductive biases [6, 13, 88]. Self-supervised models such as DINOv2 [22, 62] display greater holistic shape processing abilities [25, 37] than their supervised counterparts, implying that training may alter how information is integrated within attention. Identifying the nature and location of these changes could clarify why self-supervision yields more invariant and globally coherent representations. This motivates our third question.

Question 3. *What structural properties of intermediate-layer representations distinguish self-supervised from supervised models, and how do these properties—specifically positional preservation and content enrichment—enable holistic shape processing?*

Contributions. To address these questions, we build upon the analytical paradigm of [56, 75] (originally developed for LLMs) and contextualize it for vision. Specifically,

- **We introduce *Bi-Orthogonal Factor Decomposition (BFD)*.** A theoretical framework that couples a statistical factorization of activations with a spectral decomposition of the attention interaction matrix.
- **Quantifying information flow in attention.** BFD quantifies how much attention energy flows through positional versus content-based interactions. Content-based modes (content-content and content-position) carry the majority of energy in both architectures, confirming that contextual integration reflects semantic computation rather than positional structure. DINOv2 dedicates more energy to content-position coupling and operates through higher-rank mode spectra, suggesting richer interactions.
- **Functional specialization of attention heads.** Across models and layers, heads differentiate into content-content, content-position, and position-position operators. Within individual heads, singular modes show analogous specialization, establishing a quantitative notion of head function.
- **Dual preservation as the signature of holistic shape processing.** DINOv2 simultaneously preserves 2D spatial topology and enriches semantic content through intermediate layers. We find that the activation similarities that usually exhibit block-diagonal structure disappear when examining content alone. This reveals that content progressively enriches by integrating both positional and semantic information. Supervised models, by contrast, collapse spatial structure to quasi-1D by mid-depth, precluding this localized hierarchical integration.

To address these questions, we build upon the analytical paradigm pioneered by Song and Zhong [75] for language

models and contextualize it to ViT.

2. Theoretical Framework

We present Bi-orthogonal Factor Decomposition (BFD), which couples statistical factorization of activations with spectral decomposition of attention (Figure 2). The framework proceeds in two stages: first, we adapt ANOVA-based factorization to disentangle spatial topology from semantic content in vision representations; second, we perform SVD on the query-key interaction matrix to reveal how these factors drive token communication through bi-orthogonal modes.

Preliminaries. Let f a ViT that partition an image $x \in \mathcal{X}$ into P non-overlapping patches, each linearly projected to form a token embedding. The network admits a series of L intermediate representations

$$\mathbf{A}_\ell = f_\ell(x), \ell \in \{1, \dots, L\} \quad \text{where} \quad \mathbf{A}_\ell \in \mathbb{R}^{P \times d},$$

stacks the d -dimensional embeddings for all P tokens at layer ℓ . We denote the embedding of token p at layer ℓ as $\mathbf{A}_\ell^{(p)} \in \mathbb{R}^d$. The architecture alternates between element-wise transformations (feedforward blocks) and token-mixing operations performed by multi-head self-attention. Each block’s output is added to the residual stream, forming the standard Transformer architecture [27, 81]. Crucially, positional information must be injected into tokens; without it, attention would be permutation-invariant [86] and thus insensitive to spatial structure [20]. Most ViTs achieve this by adding a learned positional embedding $\gamma_p \in \mathbb{R}^d$ to each token at the input layer: $\mathbf{A}_0^{(p)} = \mathbf{x}^{(p)} + \gamma_p$, where $\mathbf{x}^{(p)}$ denotes the linearly projected patch embedding. This additive positional bias introduces spatial grounding and enables attention to mix tokens in a position-dependent manner. The central questions of this work concern how tokens interact through attention, and what information—positional or semantic—is exchanged during this interaction. To address this, we decompose the Transformer computation into two interpretable stages: (i) a statistical factorization of activations that disentangles position and content, and (ii) a bi-orthogonal spectral decomposition of attention that exposes its communication modes.

Statistical Factorization. Positional embeddings enter the model additively and linearly at the input layer, suggesting that subsequent activations retain a linearly recoverable positional component. We exploit this property to construct a statistically orthogonal decomposition of activations into three factors: global mean average (layer effect), positional bias (token position effect), and content residual (token idiosyncratic information). This two-ways ANOVA style decomposition provides a principled separation of signal sources [68].

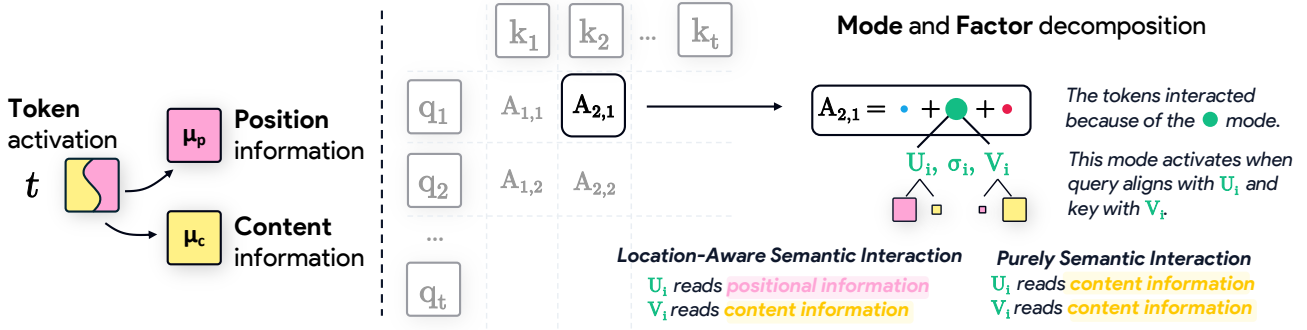


Figure 2. **Bi-orthogonal Factor Decomposition (BFD)**. Each token activation is decomposed into positional (μ_p) and content (μ_c) components. The bilinear interaction between these components is then analyzed through a singular value decomposition, yielding biorthogonal mode pairs (U_i, V_i) that explain why two tokens interact. A mode is activated when the query aligns with U_i and the key with V_i . For example, when U_i aligns with a token’s positional component while V_i aligns with another token’s content component, the mode expresses a location-aware semantic interaction, effectively mixing where something is with what it is. This decomposition exposes which paired directions enable tokens to interact and clarifies whether their interaction arises from positional or content factors or both.

Definition 1 (Vision Transformer Spatial-Content Factorization). Let \mathbb{E}_x denote expectation over images and $\mathbb{E}_{x,p}$ expectation over images and token indices. For any vision transformer f at layer ℓ , every patch embedding admits the unique additive decomposition:

$$f_\ell^{(p)}(x) = \mu_1 + \mu_p^{(p)} + \mu_c^{(p)}(x), \text{ where}$$

$$\begin{cases} \mu_1 = \mathbb{E}_{(x,p)}(f_\ell^{(p)}(x)) & (\text{Layer effect}) \\ \mu_p^{(p)} = \mathbb{E}_x(f_\ell^{(p)}(x)) - \mu_1 & (\text{Positional effect}) \\ \mu_c^{(p)}(x) = f_\ell^{(p)}(x) - \mu_1 - \mu_p^{(p)} & (\text{Content residual}) \end{cases}$$

Where the content residual $\mu_c^{(p)}(x)$ isolates input-specific semantic information and is statistically orthogonal to both global and positional components by construction. Specifically, in expectation over the joint distribution of (x, p) we have:

$$\begin{aligned} \mathbb{E}_{x,p}(\mu_1^\top \mu_p^{(p)}) &= \mathbb{E}_{x,p}(\mu_1^\top \mu_c^{(p)}(x)) \\ &= \mathbb{E}_x(\mu_p^{(p)\top} \mu_c^{(p)}(x)) = 0. \end{aligned}$$

Essentially, this factorization uses dataset and spatial marginals to isolate factors that are already implicated by the architecture. Extensions to finer-grained marginalizations, for instance, conditioning on image level activations or semantic classes are conceptually straightforward and constitute promising directions for future work.

Having decomposed activations into interpretable factors, we now require a complementary decomposition to analyze how these factors interact through attention. Specifically, we seek to decompose the query-key interaction in a manner that reveals which paired directions drive token communication at each layer and head.

Bi-orthogonal Decomposition We now analyze the mechanism that mediates information exchange: self-attention.

For a given head, the attention is then defined as

$$Y = \text{softmax}\left(\frac{(AW_Q)(AW_K)^\top}{\sqrt{d}}\right)V,$$

where we recall that $A \in \mathbb{R}^{P \times d}$ stacks token activations for P positions. Specifically, the interaction between queries and keys occurs through the bilinear form $(AW_Q)(AW_K)^\top$, which encapsulates all pairwise relations before normalization. Following the spectral analyses of [56, 63, 75], we characterize this operator via singular value decomposition:

$$(AW_Q)(AW_K)^\top = AW_QW_K^\top A^\top = AW A^\top$$

With $W \equiv W_QW_K^\top$ the interaction matrix. We decompose this matrix using SVD $W = U\Sigma V$. Each triplet (u_i, σ_i, v_i) defines a bi-orthogonal mode [24, 43]: left and right singular vectors u_i and v_i are orthonormal bases of the query and key subspaces, and σ_i quantifies their coupling strength. Crucially, the bi-orthogonality structure ensures that mode i operates independently: queries aligned with u_i communicate exclusively with keys aligned with v_i , while u_i cannot interact with v_j for $j \neq i$. This orthogonal decomposition of the attention interaction thus partitions information flow into independent (non-interfering) communication channels.

Definition 2 (Bi-orthogonal Mode Decomposition.). Given a token representation A and the interaction matrix decomposition $W \equiv W_QW_K^\top = U\Sigma V$, we define the projected codes

$$z^Q = AU, \quad z^K = AV.$$

The pair (z_i^Q, z_i^K) represents the token activations aligned with the i -th communication mode, whose strength is modulated by σ_i .

This decomposition is complete: the query-key interaction matrix can be exactly reconstructed by summing over all modes,

$$QK^\top = AWA^\top = \sum_{i=1}^d \mathbf{z}_i^Q \sigma_i(\mathbf{z}_i^K)^\top,$$

ensuring that the spectral decomposition partitions attention energy without information loss.

With both the factor decomposition of activations and the mode decomposition of attention now established, we can unify these frameworks to attribute each communication mode to its underlying informational source (position, content, or global bias), thereby exposing what information flows through each attention channel. With both the factor decomposition of activations and the mode decomposition of attention now established, we can unify these frameworks to attribute each communication mode to its underlying informational source (position, content, or global bias), thereby exposing what information flows through each attention channel.

BFD: Coupling Factors and Modes. Combining the activation factorization with the bi-orthogonal decomposition enables us to attribute each communication mode to its underlying informational factors: position, content, or global bias. Recall that each layer’s activations decompose using the three factors $\mathcal{F} = (\mu_1, \mu_p^{(p)}, \mu_c^{(p)}(\mathbf{x}))$:

$$\mathbf{A}_\ell = \mu_1 + \mu_p^{(p)} + \mu_c^{(p)}(\mathbf{x}) = \sum_{\mu \in \mathcal{F}} \mu.$$

Critically, the coupled decomposition is also complete: the query-key interaction can be exactly reconstructed by summing over all factors and modes,

$$QK^\top = \sum_{\mu \in \mathcal{F}} \sum_{i=1}^d (\mu \mathbf{u}_i) \sigma_i(\mu \mathbf{v}_i)^\top,$$

We can now project each factor of layer ℓ onto the query and key modes and measure its associated energy on the i -th mode. Formally,

$$\xi_{\cdot, i}^{(i)} = \|\mathbf{z}_{\cdot, i}^Q \cdot \sigma_i \cdot \mathbf{z}_{\cdot, i}^K\|_2^2,$$

where $\mathbf{z}_{\cdot, i}^Q$ is the projection of a factor onto query mode \mathbf{u}_i (e.g., $\mu_p^{(p)} \mathbf{U}_i$ for position), similarly for keys. The normalized energy:

$$\bar{\xi}_{\cdot}^{(i)} = \mathbb{E}_{\mathbf{x}} \left(\frac{\xi_{\cdot}^{(i)}(\mathbf{x})}{\sum_j \xi_{\cdot}^{(j)}(\mathbf{x})} \right), \quad \sum_i \bar{\xi}_{\cdot}^{(i)} = 1.$$

quantifies the relative dominance of each factor within mode i (e.g., $\bar{\xi}_{\cdot}^{(i)}$ is the relative energy of the content information for the i -th mode). Aggregating across modes and heads yields layer-level summaries that characterize the informational profile of attention at different depths. These quantities operationalize statements such as "this head is content-dominated" or "this layer pivots from position to content," providing a quantitative backbone linking representational factorization to emergent specialization in Vision Transformers.

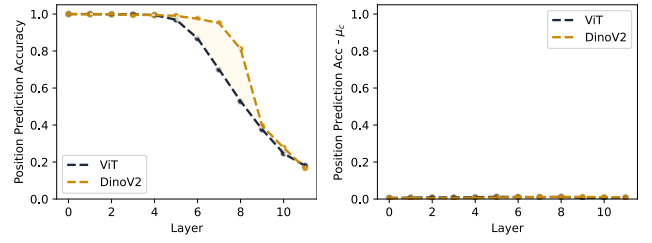


Figure 3. Validating Factor Isolation. Linear probes decode spatial coordinates from factorized representations across layers. *Left:* Position is successfully decoded from the raw block activations in both architectures, with DINOv2 maintaining higher accessibility through intermediate layers before both converge to 20% accuracy at depth. *Right:* Content factors $\mu_c^{(p)}(\mathbf{x})$ contain negligible positional information (chance-level decoding), confirming that the factorization removes positional information and achieves statistical orthogonality. The asymmetric preservation trajectory reveals that self-supervised training maintains linearly accessible positional structure deeper into the network than supervised learning.

3. Results

Experimental Setup. We apply our framework to two representative Vision Transformers: a supervised ViT-B/16 from timm library [85] and a self-supervised DINOv2-B/14 with registers [22]. Both models share identical architectural parameters: 12 transformer blocks, 16×16 patch tokenization in the supervised ViT-B/16 and 14×14 patch tokenization yielding 196 and 256 spatial tokens, plus a class token (and 4 additional register tokens for DINOv2). We compute the statistical factorization (Definition 1) and bi-orthogonal mode decomposition (Definition 2) across all 12 layers for both architectures, analyzing 5,000 images sampled from the ImageNet validation set.

Sanity Check: Validating the Decomposition. Before interrogating attention mechanisms, we first validate that our ANOVA-based factorization genuinely isolates position from content. Figure 3 presents two complementary tests. First, we train linear probes to decode spatial coordinates directly from the activations (\mathbf{A}): both ViT and DINOv2 exhibit near-perfect position availability in early layers (left panel), confirming that spatial coordinates are linearly recoverable from the activations themselves. However, a key difference emerges—DINOv2 maintains substantially higher decoding accuracy through intermediate layers (layers 5-8), while ViT’s positional availability degrades more rapidly. Both architectures eventually converge to approximately 20% accuracy in the deepest layers, suggesting that extreme depth compresses positional information regardless of training paradigm. Critically, the content residual $\mu_c^{(p)}(\mathbf{x})$ contains negligible positional information across all layers in both models – probes trained on content factors alone achieve chance-level decoding accuracy (right panel, performance equivalent to random guessing). This orthogonality by construction, validated empirically, establishes that our decomposition cleanly disentangles the two factors without

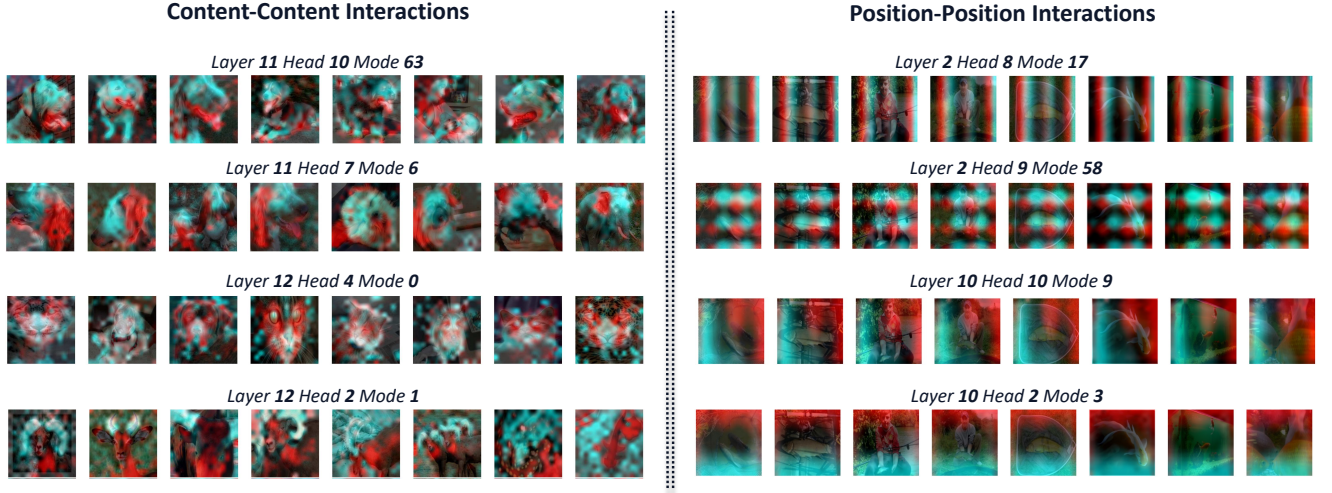


Figure 4. **Projections of bi-orthogonal modes in DINOv2.** The query (red) and key (cyan) singular vectors are projected onto either the content or positional factor, highlighting the image regions that most strongly activate each singular direction.

information leakage. The sustained linear accessibility of position from the block activations in DINOv2’s intermediate layers provides a first hint that self-supervised learning may preserve spatial structure more *explicitly* than supervised training, a hypothesis we will examine geometrically in our analysis of [Question 3](#). Having confirmed proper factor isolation, we now examine how these factors participate in attention-mediated communication.

Characterizing Information Flow. [Question 1](#) asks which informational factors—layer effect, position, or content—mediate token communication through attention. We address this by projecting each mode’s query and key directions $(\mathbf{u}_i, \mathbf{v}_i)$ onto the three factors and compute the normalized energy contribution of each bi-orthogonal mode for all unique interactions (see [Appendix C](#)). We then aggregate across modes and heads to obtain layer-level profiles ([Figure 5](#)). Both architectures allocate the majority of attention energy to content-based interactions: content-content (pure

semantic exchange) and content-position (localization-aware semantic integration) together dominate the interaction budget, particularly in deeper layers. This finding confirms that contextual integration in later blocks reflects genuine content-driven computation rather than mere positional structure – the content factor $\mu_c^{(p)}(\mathbf{x})$, which varies across images, accounts for the bulk of query-key alignment (for qualitative examples of these interactions, see [Figure 1](#), [Figure 4](#) and [Appendix Sec. A](#)). However, a divergence emerges between training paradigms: DINOv2 dedicates greater energy to content-position than ViT across all layers. At a finer resolution, we observe that this energy is distributed broadly across many modes rather than concentrated in a few dominant ones. As shown in [Figure 6](#), DINOv2 spreads its content-based interaction energy across a wide range of singular modes, indicating that heads rely on a diverse set of interaction patterns rather than a small set of strong modes. In contrast, the supervised ViT not only allocates less energy to content-based interactions overall, but this energy is also more concentrated in a few dominant modes (see [Appendix Figure 17](#)).

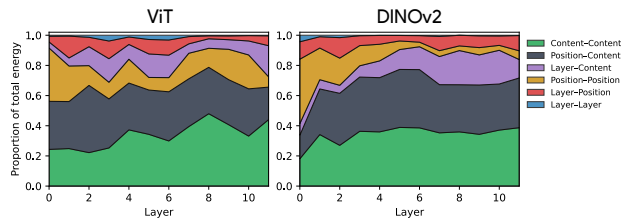


Figure 5. **Energy Distribution Across Informational Factors.** Layer-wise decomposition of attention energy into contributions from layer effect (global bias), position, content, and their pairwise interactions. Both architectures allocate the majority of energy to content-based interactions (content-content and content-position), confirming that attention-mediated interaction reflects genuine semantic computation rather than positional structure alone. A divergence emerges: DINOv2 dedicates substantially greater energy to content-position coupling across all layers, which could be a sign of localization enrichment of semantic content (content gets modulated by positional information).

To understand whether this energy distribution reflects distributed computation or modular specialization, we next examine the spectral structure of the interaction matrix $\mathbf{W} = \mathbf{W}_Q \mathbf{W}_K^\top$ itself. [Figure 7](#) (left panels) reveals that DINOv2 maintains substantially higher stable rank than ViT across layers, indicating that its attention operates through a richer, more distributed set of communication modes rather than concentrating interaction on a few dominant directions. This higher rank is a direct consequence of a flatter singular value spectrum; DINOv2 sustains a greater number of modes with high singular values compared to the supervised ViT, whose spectrum decays more rapidly (see [Appendix Figure 21](#) for a detailed layer-by-layer view). The alignment between query and key modes, quantified as $\cos(\mathbf{u}_i, \mathbf{v}_i) \cdot \sigma_i$, further distinguishes the two models ([Figure 7](#); right panel):

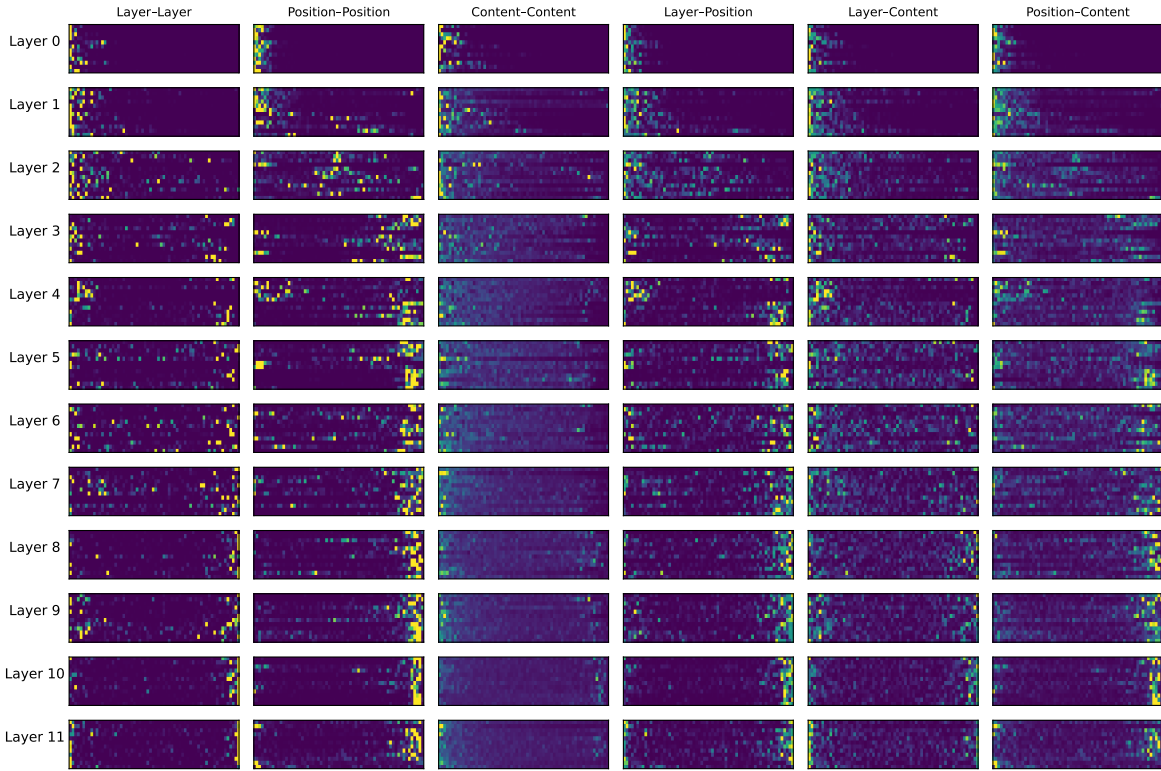


Figure 6. **Rich Content-Driven Interactions in DINOv2.** Each subplot visualizes the normalized energy for a specific interaction factor (columns) at each layer (rows) for all the decomposed modes. For each layer and interaction, the x-axis represents the bi-orthogonal modes, sorted by singular value from highest (left) to lowest (right) and the y axis represents the head. The color intensity shows the relative contribution of a single mode to a given interaction’s total energy (normalized horizontally for each head). The visualization reveals that the low-singular-value modes (right side) also become increasingly dedicated to carrying content-based information in deeper layers.

ViT exhibits high alignment (modes are nearly symmetric), whereas DINOv2 displays lower alignment, particularly in intermediate layers, signaling asymmetric query-key specialization. This spectral signature: higher rank, asymmetric modes, and high content-position energy, constitutes a quantitative fingerprint of self-supervised attention. Yet the existence of multiple modes raises a natural question: do these modes serve distinct functional roles, or do they redundantly encode the same information?

Functional Specialization: Mode Purity and Head Differentiation. **Question 2** concerns whether attention heads and their constituent singular modes exhibit functional specialization into distinct informational operators. We address this by projecting each mode’s query and key directions (u_i, v_i) onto the three informational factors – layer, position, and content and computing the normalized interaction energies for the six undirected factor pairs (see Appendix Sec. C). We then group these six interactions into three families (layer, position, and content) and compute barycentric coordinates that indicate the relative contribution of layer-, position-, and content-related interactions for each mode. Visualizing the distribution of modes in these coordinates using ternary plots (Figure 8) — where the vertices correspond to layer-dominated, position-dominated, and content-

dominated operators, and points along edges reflect mixed operators (e.g., content–position)—reveals strong mode purity: individual modes cluster near specific vertices rather than dispersing uniformly across the simplex.

This indicates that most modes implement specialized informational operations (Appendix Sec. A): some mediate pure semantic exchange (content-content), others mainly track spatial relationships (position-position), and still others integrate semantics with spatial context (content-with-position). Crucially, this specialization holds both at the mode level (individual singular vectors within a head) and at the head level (when aggregating energy across all modes of a head). Across both ViT and DINOv2, heads differentiate into three functional classes, with DINOv2 exhibiting slightly greater mode purity and more pronounced differentiation in intermediate layers (Appendix Figure 18 and Appendix Figure 20 provide per-layer plots). This functional modularity implies that per-head mechanistic analyses are not fundamentally incomplete; BFD decomposition of heads finds pure computational primitives that can be meaningfully composed. The consistency of specialization across architectures suggests it is an emergent property of multi-head self-attention rather than a training-specific artifact, though the degree of purity varies with supervision. Having established that attention mechanisms decompose into

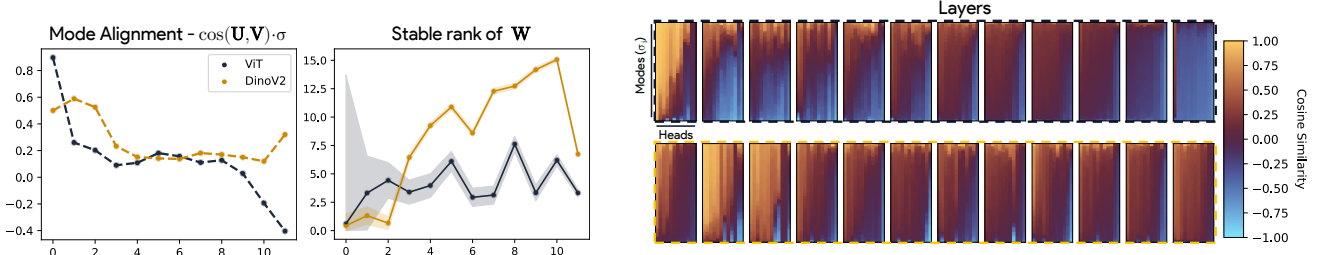


Figure 7. **Spectral Structure of Query-Key Interactions.** *Left:* Mode alignment $\cos(\mathbf{u}_i, \mathbf{v}_i) \cdot \sigma_i$ across layers. Both architectures transition from high alignment (symmetric modes) in early layers to near-orthogonality at depth, consistent with prior observations [63]. A subtle but important difference emerges: ViT exhibits negative cosine values (anti-aligned modes), while DINOv2 approaches contextual grouping behavior with modes clustering near orthogonality rather than anti-alignment. *Middle:* Stable rank of the interaction matrix \mathbf{W} reveals that DINOv2 maintains substantially higher effective dimensionality across all layers, indicating that its attention operates through a richer subspace with more distributed communication channels. *Right:* Mode-specific alignment distributions refine the aggregate view. ViT concentrates mass in negative-alignment modes, whereas DINOv2 exhibits a more balanced distribution. Critically, within each architecture, modes differentiate into specialized regimes – some highly aligned (symmetric query-key relationships), others orthogonal or anti-aligned – establishing mode-level functional diversity beyond head-level aggregates.

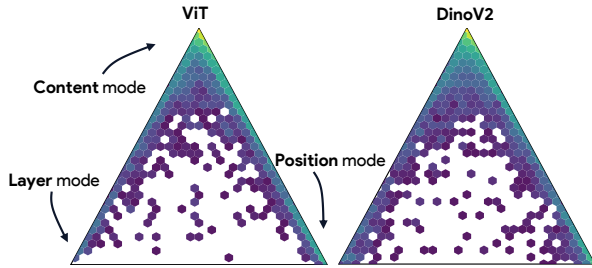


Figure 8. **Functional specialization of Attention Modes.** Ternary plots showing the barycentric coordinates of each bi-orthogonal mode across all the heads and layers in ViT (left) and DINOv2 (right). Each point reflects the relative contribution of layer-, position-, and content-related interactions for a mode. Modes cluster near the vertices rather than dispersing uniformly across the simplex, revealing strong functional specialization for both the models, with tighter clustering in DINOv2, indicating more pronounced specialization under self-supervised training.

functionally interpretable operators, we now ask: what representational properties enable DINOv2’s superior holistic shape processing, and how do these properties manifest in the geometry of learned representations?

Dual Preservation: Positional Structure and Semantic Enrichment as the Signature of Holistic Processing.

Question 3 asks what structural properties of intermediate-layer representations distinguish self-supervised from supervised models and enable holistic shape processing [25, 37]. To probe the geometry of positional representations, we apply PCA to the positional factor $\mu_p^{(p)}$ across all 196 and 256 spatial tokens and visualize the first three principal components as RGB-encoded 3D point clouds (Figure 9). A divergence emerges: DINOv2 preserves the 2D spatial grid topology throughout all layers—tokens retain their neighborhood structure, and the manifold remains approximately planar (2D sheet embedded in 3D space) even at depth. In

contrast, ViT progressively collapses the spatial manifold into a quasi-1D structure by layer 5, with tokens condensing along a single principal axis and obliterating local spatial relationships (Appendix Figure 22 and Figure 23 provide rotated views across all layers). While linear probes confirm that positional information remains decodable in both models, the collapse of the positional subspace dimensionality in ViT suggests its spatial coordinate system becomes far less expressive in mid-to-later layers. This positional preservation in DINOv2 directly addresses the first half of Question 3: self-supervised models maintain “explicit” spatial structure where supervised models do not. Yet preservation alone does not explain holistic processing—the model must also enrich semantic content.

Next, we examine how semantic content evolves across the network. We analyze this by computing the Pearson correlation of content factors and the full activations, of all tokens, for each pair of layers ($r(\mu_c^{(p)}(\mathbf{x})_i, \mu_c^{(p)}(\mathbf{x})_j)$ and $r(\mathbf{A}_i, \mathbf{A}_j)$) in Figure 10. While the full activations show a higher degree of similarity across all layers, with similarity decaying slowly from the diagonal, the similarity matrix for the content factor exhibits a smoother enrichment trajectory, with incremental refinement through the intermediate layers. The correlation between early- and late-layer’s content factor is significantly lower than for the full activations. This indicates that much of the apparent similarity in the full activations is an artifact of the persistent positional code. Once the positional signal is removed, it is clear that the underlying semantic content is undergoing a transformation as it progresses through each layer the network.

In contrast, the ViT similarity structure reveals a different pattern. The layer–layer correlations for the content factors closely mirror those of the full activations, despite the content factors being position-free. This indicates that positional information plays a far smaller role in shaping ViT’s internal geometry, consistent with the early collapse of the positional scaffold we observed earlier. More importantly, the con-

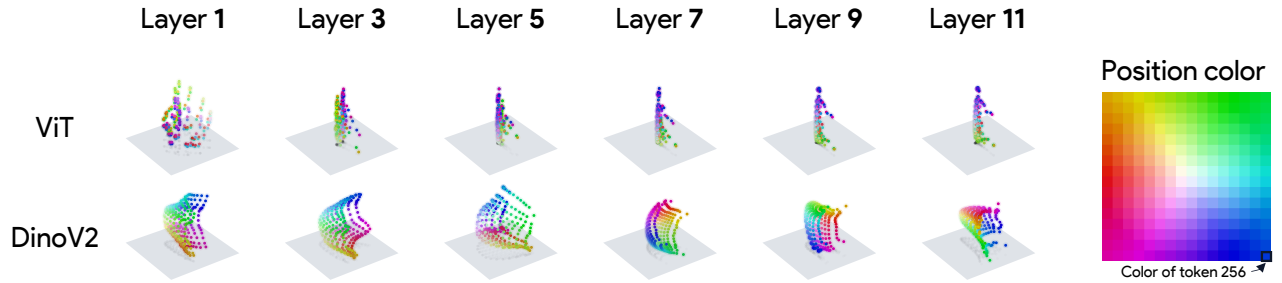


Figure 9. **Geometric Preservation of Positional Structure.** PCA visualization of positional factors $\mu_p^{(p)}$ across layers, with each token’s spatial location encoded by color. The first three principal components are rendered as RGB-encoded 3D point clouds. *Top (ViT)*: Supervised training progressively collapses the spatial manifold—by layer 5, the 2D grid degenerates into a quasi-1D structure with tokens condensing along a single principal axis, obliterating neighborhood relationships. *Bottom (DINOv2)*: Self-supervised training preserves the 2D spatial grid topology throughout all layers. Tokens retain their neighborhood structure, and the manifold remains approximately planar (a 2D sheet embedded in 3D space) even at depth.

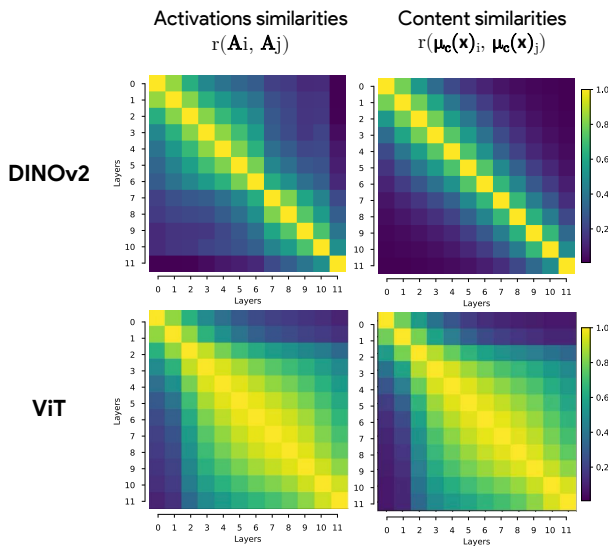


Figure 10. **Content Enrichment with Depth.** Pairwise Pearson correlations across layers for full activations (A_i) and content factors ($\mu_c(x)_i$). *DINOv2 (top)* shows a gradual transformation of content across depth, revealed once positional signals are removed. In *ViT (bottom)*, content and activations show similar trajectories, indicating weaker positional influence and limited semantic evolution.

tent factors themselves show only limited evolution across depth: early and late layers remain highly correlated. This suggests that, rather than progressively enriching semantic content through contextual interactions, the ViT primarily performs more local refinement (e.g., refining an “eye” into a slightly sharper “eye”) rather than integrating it into a broader semantic configuration (e.g., placing the eye within the context of a face). DINOv2, by contrast, shows a steady transformation of its content factors across layers, indicative of contextually enriched representations: token meaning is updated by accumulating information from other parts of the object (see example for semantic interactions in Appendix Figure 13). This rapid evolution of semantic content in DI-

NOv2, is the second key component of the dual preservation. We believe that the dual preservation – maintaining a stable spatial topology while enabling the evolution of semantic content – constitutes the computational motif underlying holistic shape processing.

4. Discussion

We introduced Bi-orthogonal Factor Decomposition (BFD), a framework that jointly factorizes vision transformer activations into positional and content components and spectrally decomposes their attention interactions. This coupling exposes what type of information flows through self-attention and how it is organized across heads and modes. When applied to supervised and self-supervised ViTs, BFD reveals a consistent computational structure within self-attention. First, heads and their singular modes exhibit clear functional specialization: rather than mixing informational factors, each channel reliably operates as a position, content, or mixed content–position operator. Second, both architectures allocate most of their attention energy to content–dominated interactions, indicating that semantic features, not positional offsets, drive the majority of token-to-token communication. Third, this pattern is amplified in DINOv2: content-dominated interactions are stronger and distributed across a larger set of modes, reflecting a richer and more distributed computational motif than in supervised ViTs. Finally, self-supervised models maintain a well-structured positional scaffold while their content representations progressively enrich across depth, enabling coordinated spatio-semantic integration that is largely absent in supervised counterparts, hence offering a mechanistic explanation for their superior holistic shape processing.

References

- [1] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. *Advances in Neural Information Processing Systems (NIPS)*, 2018.

- [2] Arjun R Akula and Song-Chun Zhu. Attention cannot be an explanation. *ArXiv e-print*, 2022.
- [3] Ibrahim M Alabdulmohsin, Xiaohua Zhai, Alexander Kolesnikov, and Lucas Beyer. Getting vit in shape: Scaling laws for compute-optimal model design. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [4] Lucian Alecu, Hugues Bonnin, Thomas Fel, Laurent Gardes, Sébastien Gerchinovitz, Ludovic Ponsolle, Franck Mamalet, Éric Jenn, Vincent Mussot, Cyril Cappi, et al. Can we reconcile safety objectives with machine learning performances? *ERTS*, 2022.
- [5] Xiaoyan Bai, Itamar Pres, Yuntian Deng, Chenhao Tan, Stuart Shieber, Fernanda Viégas, Martin Wattenberg, and Andrew Lee. Why can’t transformers learn multiplication? reverse-engineering reveals long-range dependency pitfalls. *ArXiv e-print*, 2025.
- [6] Adrien Bardes, Jean Ponce, and Yann LeCun. Vircreg: Variance-invariance-covariance regularization for self-supervised learning. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2022.
- [7] Jasmijn Bastings and Katja Filippova. The elephant in the interpretability room: Why use attention as explanation when we have saliency methods? *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, 2020.
- [8] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [9] Adrien Bibal, Rémi Cardon, David Alfter, Rodrigo Wilkens, Xiaoou Wang, Thomas Francois, and Patrick Watrin. Is attention explanation? an introduction to the debate. *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2022.
- [10] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermy, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023.
- [11] Lucius Bushnaq, Dan Braun, and Lee Sharkey. Stochastic parameter decomposition. *ArXiv e-print*, 2025.
- [12] Charles F Cadieu, Ha Hong, Daniel LK Yamins, Nicolas Pinto, Diego Ardila, Ethan A Solomon, Najib J Majaj, and James J DiCarlo. Deep neural networks rival the representation of primate it cortex for core visual object recognition. *PLoS computational biology*, 2014.
- [13] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [14] Edward Choi, Mohammad Taha Bahadori, Jimeng Sun, Joshua Kulas, Andy Schuetz, and Walter Stewart. Retain: An interpretable predictive model for healthcare using reverse time attention mechanism. *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [15] Brianna Chrisman, Lucius Bushnaq, and Lee Sharkey. Identifying sparsely active circuits through local loss landscape decomposition. *ArXiv e-print*, 2025.
- [16] Michael Chui, Eric Hazan, Roger Roberts, Alex Singla, and Kate Smaje. The economic potential of generative ai. 2023.
- [17] SueYeon Chung. Neural population geometry: An approach for understanding biological and artificial neural networks. *Current opinion in neurobiology*, 2021.
- [18] Julien Colin, Thomas Fel, Rémi Cadène, and Thomas Serre. What i cannot predict, i do not understand: A human-centered evaluation framework for explainability methods. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [19] Colin Conwell, Jacob S Prince, Kendrick N Kay, George A Alvarez, and Talia Konkle. What can 1.8 billion regressions tell us about the pressures shaping high-level visual representation in brains and machines? *Nature communications*, 2022.
- [20] Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. On the relationship between self-attention and convolutional layers. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [21] Valérie Costa, Thomas Fel, Ekdeep Singh Lubana, Bahareh Tolooshams, and Demba Ba. From flat to hierarchical: Extracting sparse representations with matching pursuit. *Advances in Neural Information Processing Systems (NeurIPS)*, 2025.
- [22] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. *ArXiv e-print*, 2023.
- [23] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. *Proceedings of the International Conference on Machine Learning (ICML)*, 2023.
- [24] Jean Dieudonné. On biorthogonal systems. *Michigan Mathematical Journal*, 1953.
- [25] Fenil R Doshi, Thomas Fel, Talia Konkle, and George Alvarez. Visual anagrams reveal hidden differences in holistic shape processing across vision models. *Advances in Neural Information Processing Systems (NeurIPS)*, 2025.
- [26] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *ArXiv e-print*, 2017.
- [27] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [28] Maximilian Dreyer, Jim Berend, Tobias Labarta, Johanna Vielhaben, Thomas Wiegand, Sebastian Lapuschkin, and Wojciech Samek. Mechanistic understanding and validation of large ai models with semanticlens. *ArXiv e-print*, 2025.

- [29] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022.
- [30] Thomas Fel, Remi Cadene, Mathieu Chalvidal, Matthieu Cord, David Vigouroux, and Thomas Serre. Look at the variance! efficient black-box explanations with sobol-based sensitivity analysis. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [31] Thomas Fel, Ivan Felipe, Drew Linsley, and Thomas Serre. Harmonizing the object recognition strategies of deep neural networks with humans. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [32] Thomas Fel, Thibaut Boissin, Victor Boutin, Agustin Picard, Paul Novello, Julien Colin, Drew Linsley, Tom Rousseau, Rémi Cadène, Laurent Gardes, and Thomas Serre. Unlocking feature visualization for deeper networks with magnitude constrained optimization. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [33] Thomas Fel, Victor Boutin, Mazda Moayeri, Remi Cadene, Louis Bethune, Mathieu Chalvidal, and Thomas Serre. A holistic approach to unifying automatic concept extraction and concept importance estimation. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [34] Thomas Fel, Agustin Picard, Louis Bethune, Thibaut Boissin, David Vigouroux, Julien Colin, Rémi Cadène, and Thomas Serre. Craft: Concept recursive activation factorization for explainability. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [35] Thomas Fel, Louis Bethune, Andrew Kyle Lampinen, Thomas Serre, and Katherine Hermann. Understanding visual feature reliance through the lens of complexity. *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [36] Ruth C. Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [37] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [38] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2017.
- [39] Leilani H. Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. Explaining explanations: An overview of interpretability of machine learning. *Proceedings of the IEEE International Conference on data science and advanced analytics (DSAA)*, 2018.
- [40] Chris Hamblin, Thomas Fel, Srijani Saha, Talia Konkle, and George Alvarez. Feature accentuation: Revealing ‘what’ features respond to in natural images. *ArXiv e-print*, 2024.
- [41] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [42] Alon Jacovi and Yoav Goldberg. Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness? *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020.
- [43] Samuel Kaplan. Biorthogonality and integration. *Proceedings of the American Mathematical Society*, 1956.
- [44] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.
- [45] Sunnie S. Y. Kim, Nicole Meister, Vikram V. Ramaswamy, Ruth Fong, and Olga Russakovsky. HIVE: Evaluating the human interpretability of visual explanations. *Proceedings of the IEEE European Conference on Computer Vision (ECCV)*, 2022.
- [46] Talia Konkle and George A Alvarez. A self-supervised domain-general learning framework for human ventral stream representation. *Nature communications*, 2022.
- [47] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.
- [48] Matthew Kowal, Achal Dave, Rares Ambrus, Adrien Gaidon, Konstantinos G Derpanis, and Pavel Tokmakov. Understanding video transformers via universal concept discovery. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [49] Matthew Kowal, Richard P Wildes, and Konstantinos G Derpanis. Visual concept connectome (vcc): Open world concept discovery and their interlayer connections in deep models. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [50] Nikolaus Kriegeskorte. Deep neural networks: a new framework for modeling biological vision and brain information processing. *Annual review of vision science*, 2015.
- [51] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 2015.
- [52] Michael Lepori, Thomas Serre, and Ellie Pavlick. Break it down: Evidence for structural compositionality in neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [53] Michael A Lepori, Ellie Pavlick, and Thomas Serre. Neurosurgeon: A toolkit for subnetwork analysis. *ArXiv e-print*, 2023.
- [54] Drew Linsley, Ivan F Rodriguez Rodriguez, Thomas Fel, Michael Arcaro, Saloni Sharma, Margaret Livingstone, and Thomas Serre. Performance-optimized deep neural networks are evolving into worse models of inferotemporal visual cortex. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [55] Callum McDougall, Arthur Conmy, Cody Rushing, Thomas McGrath, and Neel Nanda. Copy suppression: Comprehensive understanding an attention head. *ArXiv e-print*, 2023.

- [56] Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. Talking heads: Understanding inter-layer communication in transformer language models. *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [57] Lukas Muttenthaler, Jonas Dippel, Lorenz Linhardt, Robert A Vandermeulen, and Simon Kornblith. Human alignment of neural network representations. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023.
- [58] Sabine Muzellec, Leo Andeol, Thomas Fel, Rufin VanRullen, and Thomas Serre. Gradient strikes back: How filtering out high frequencies improves explanations. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.
- [59] Giang Nguyen, Daeyoung Kim, and Anh Nguyen. The effectiveness of feature attribution methods and its correlation with automatic evaluation scores. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [60] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2017.
- [61] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *ArXiv e-print*, 2022.
- [62] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *ArXiv e-print*, 2023.
- [63] Xu Pan, Aaron Philip, Ziqian Xie, and Odelia Schwartz. Dissecting query-key interaction in vision transformers. *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [64] Jayneel Parekh, Pegah Khayatan, Mustafa Shukor, Alasdair Newson, and Matthieu Cord. A concept-based explainability framework for large multimodal models. *ArXiv e-print*, 2024.
- [65] Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. *Proceedings of the International Conference on Machine Learning (ICML)*, 2024.
- [66] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. *Proceedings of the British Machine Vision Conference (BMVC)*, 2018.
- [67] Jürgen Schmidhuber. The speed prior: a new simplicity measure yielding near-optimal computable predictions. *International conference on computational learning theory*, 2002.
- [68] George AF Seber. *Multivariate observations*. John Wiley & Sons, 2009.
- [69] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Gradcam: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [70] Thomas Serre. *Learning a dictionary of shape-components in visual cortex: Comparison with neurons, humans and machines*. PhD thesis, Massachusetts Institute of Technology, 2006.
- [71] Mustafa Shukor and Matthieu Cord. Implicit multimodal alignment: On the generalization of frozen llms to multimodal inputs. *ArXiv e-print*, 2024.
- [72] Leon Sixt, Maximilian Granz, and Tim Landgraf. When explanations lie: Why many modified bp attributions fail. *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- [73] Leon Sixt, Martin Schuessler, Oana-Iuliana Popescu, Philipp Weiß, and Tim Landgraf. Do users benefit from interpretable vision? a user study, baseline, and dataset. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2022.
- [74] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- [75] Jiajun Song and Yiqiao Zhong. Uncovering hidden geometry in transformers via disentangling position and context. *ArXiv e-print*, 2023.
- [76] Ilia Sucholutsky, Lukas Muttenthaler, Adrian Weller, Andi Peng, Andreea Bobu, Been Kim, Bradley C Love, Erin Grant, Jascha Achterberg, Joshua B Tenenbaum, et al. Getting aligned on representational alignment. *ArXiv e-print*, 2023.
- [77] Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. *Proceedings of the IEEE European Conference on Computer Vision (ECCV)*, 2022.
- [78] Martin Tutek and Jan Snajder. Staying true to your word:(how) can attention become explanation? *ArXiv e-print*, 2020.
- [79] Shimon Ullman, Liav Assif, Ethan Fetaya, and Daniel Harari. Atoms of recognition in human and computer vision. *Proceedings of the National Academy of Sciences*, 2016.
- [80] Lucrezia Valeriani, Diego Doimo, Francesca Cuturello, Alessandro Laio, Alessio Ansuini, and Alberto Cazzaniga. The geometry of hidden representations of large transformer models. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [81] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [82] Johanna Vielhaben, Stefan Blücher, and Nils Strodthoff. Multi-dimensional concept discovery (mcd): A unifying framework with completeness guarantees. *The Journal of Transactions on Machine Learning Research (TMLR)*, 2023.
- [83] Yue Wang and Sai Ho Chung. Artificial intelligence in safety-critical systems: a systematic review. *Industrial Management & Data Systems*, 2022.
- [84] Sarah Wiegrefe and Yuval Pinter. Attention is not not explanation. *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019.
- [85] Ross Wightman et al. Pytorch image models, 2019.
- [86] Hengyuan Xu, Liyao Xiang, Hangyu Ye, Dixi Yao, Pengzhi Chu, and Baochun Li. Permutation equivariance of transformers and its applications. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

- [87] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. *Proceedings of the International Conference on Machine Learning (ICML)*, 2015.
- [88] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- [89] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. *Proceedings of the IEEE European Conference on Computer Vision (ECCV)*, 2014.
- [90] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [91] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2023.
- [92] Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. A tale of two features: Stable diffusion complements dino for zero-shot semantic correspondence. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [93] Ruihan Zhang, Prashan Madumal, Tim Miller, Krista A Ehinger, and Benjamin IP Rubinstein. Invertible concept-based explanations for cnn models with non-negative concept activation vectors. *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2021.
- [94] Boyang Zheng, Nanye Ma, Shengbang Tong, and Saining Xie. Diffusion transformers with representation autoencoders. *ArXiv e-print*, 2025.

A. Qualitative Visualization of Interactions from Different Information Factors

To complement the quantitative analyses in the main paper, we present qualitative visualizations of three key interactions extracted by BFD: content–content, content–position, and position–position. For each type, we show representative bi-orthogonal modes whose query–key directions (u_i, v_i) are activated by that specific interaction. Each visualization overlays the mode’s spatial pattern onto the activating images, providing intuition about what kind of signal the mode reads from the query and the key token.

Content–Content (Purely Semantic Interactions).

These modes activate on meaningful interactions between object parts — faces, limbs, border edges — indicating semantic affinity between regions. Because both query and key directions respond to content structure, they highlight *what* the object is, not *where* it is. Qualitatively, they reveal part-to-part or part-to-whole correspondences (e.g., eye-to-head or ear-to-snout), demonstrating how semantic information is exchanged between tokens.

Position–Position (Purely Spatial Interactions). These modes activate on geometric structure rather than semantics. They highlight global spatial patterns such as left-to-right sweeps, top-to-bottom gradients, or Fourier-like waveforms. These interactions allow the modes to keep a track of how positional information flows through tokens.

Content–Position (Localization-Aware Semantic Interactions). These modes activate on semantic regions in a way that depends on spatial context. One side (query or key) aligns with semantic content while the other aligns with positional structure. As a result, they highlight meaningful regions that are modulated by where they occur in the image — e.g., activating on specific object parts based on stereo-typed spatial queries. These patterns illustrate how semantic features directly interact with positional cues.

Together, these qualitative examples show how different interaction types manifest at the mode level—semantic exchange, spatial exchange, and spatially conditioned semantic exchange—and provide intuition for the functional specialization results reported in the main text.

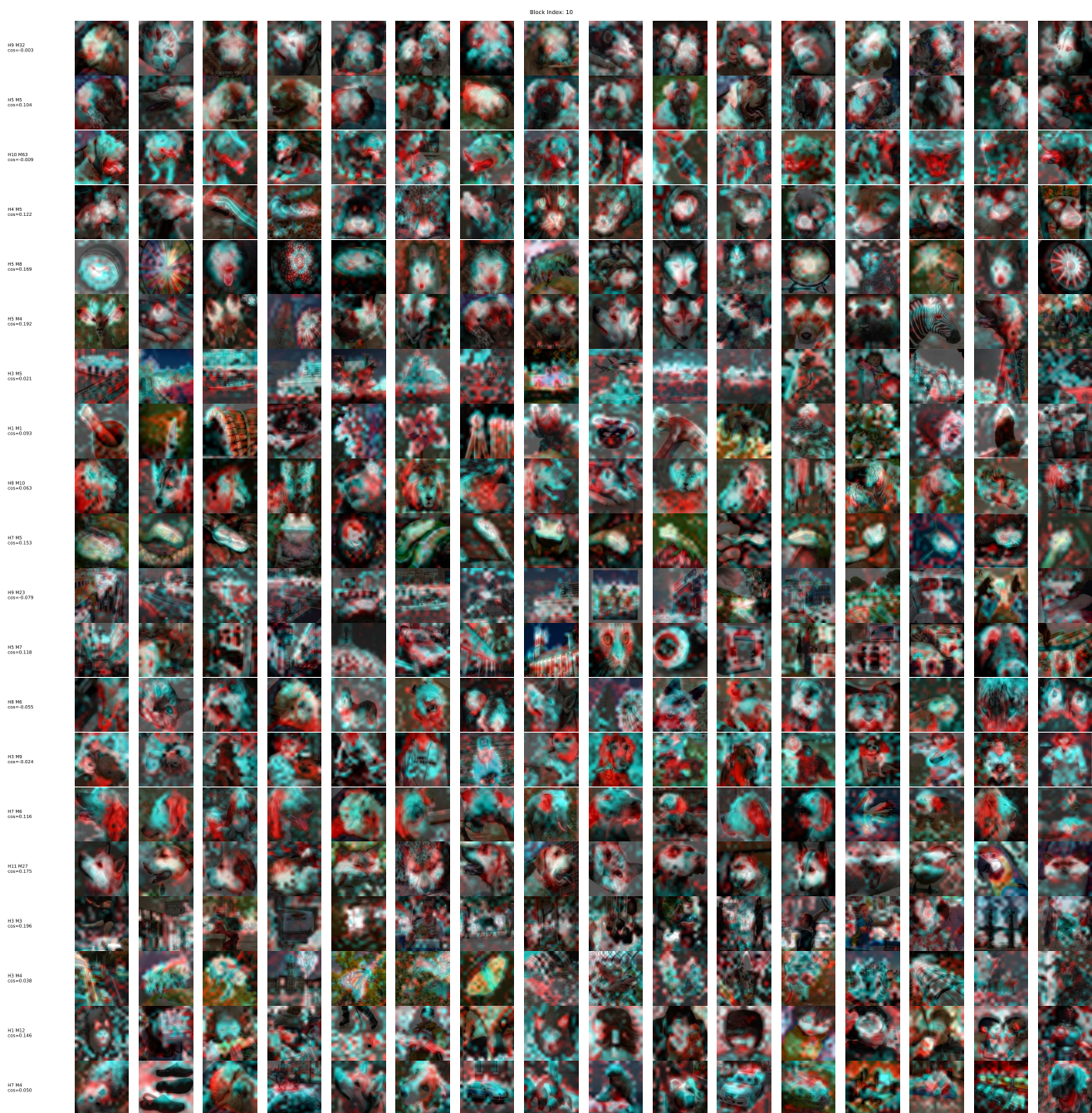


Figure 11. Content-Content Interactions in DINOv2 Block Index 10.



Figure 12. Content-Content Interactions in DINOv2 Block Index 11.

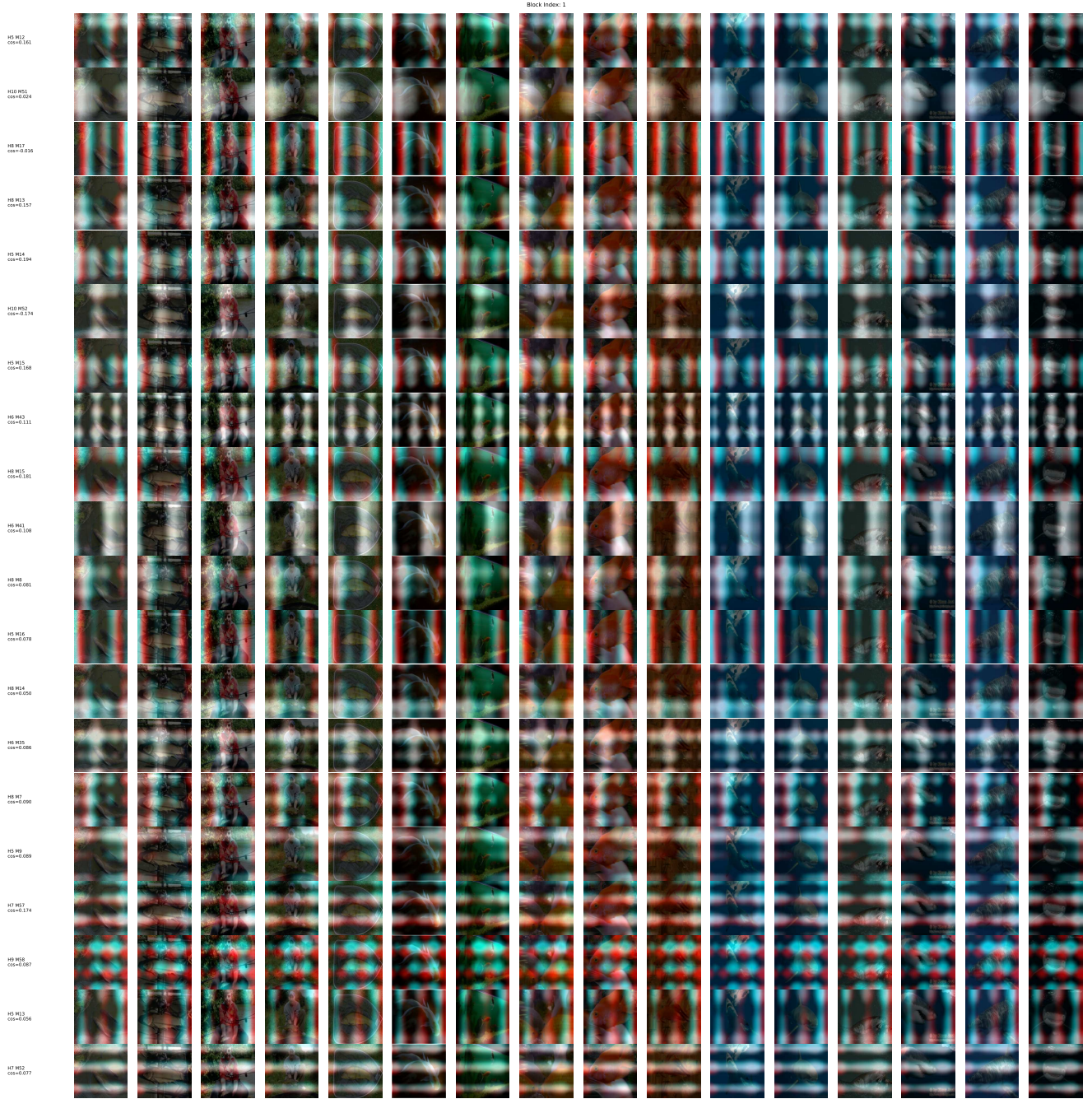


Figure 13. Position-Position Interactions in DINOv2 Block Index 1.

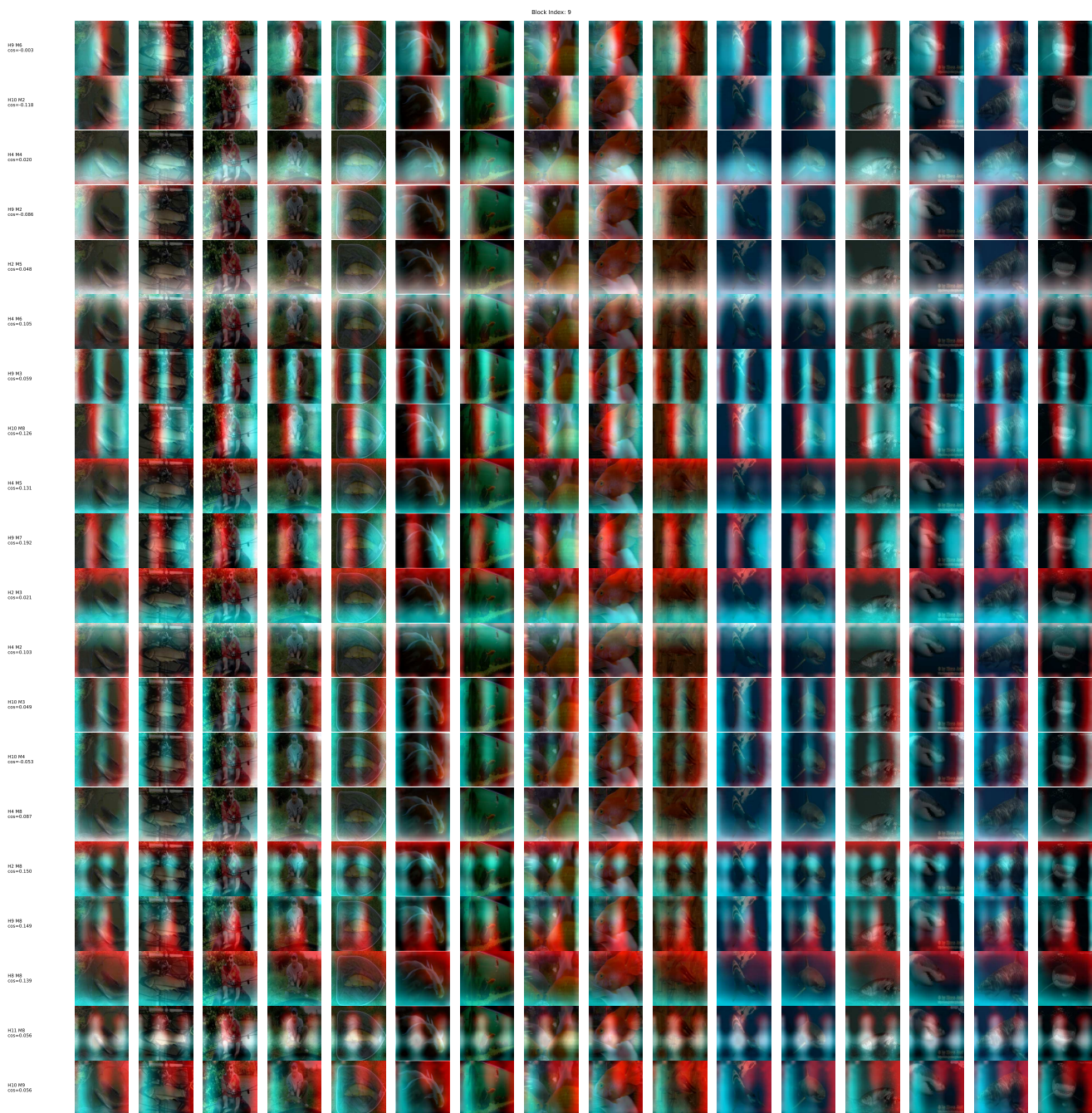


Figure 14. Position-Position Interactions in DINOv2 Block Index 9.

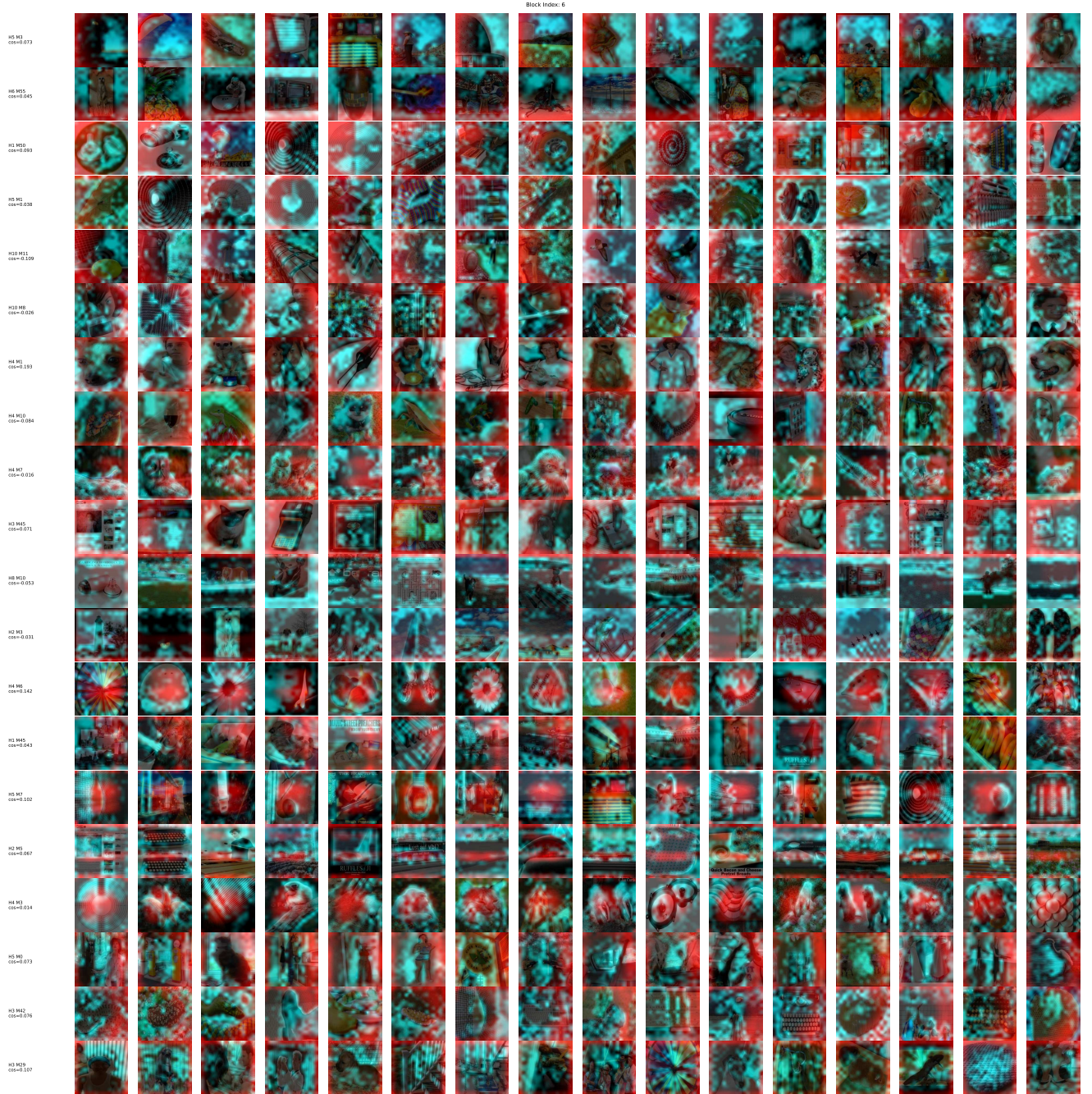
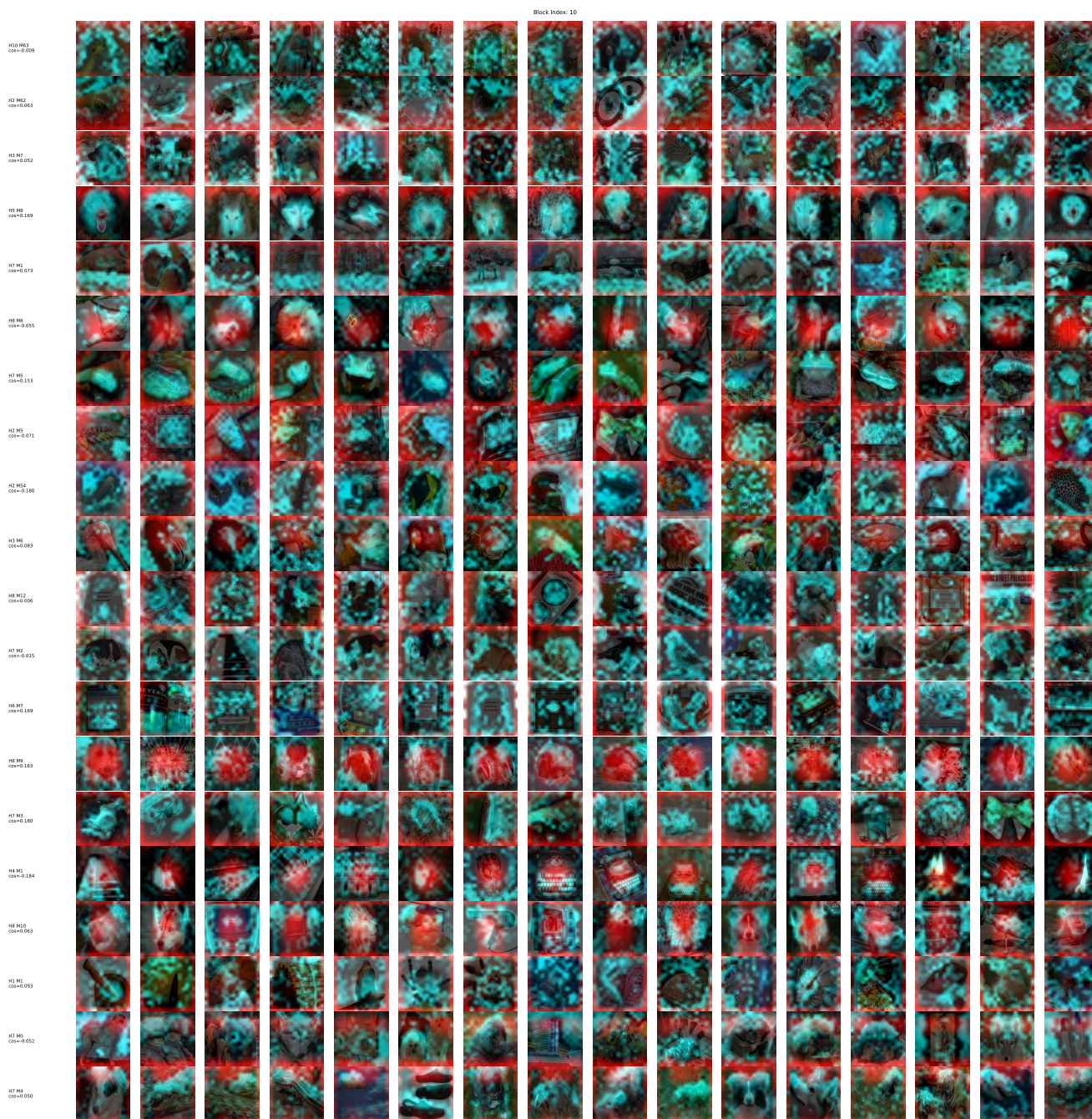


Figure 15. Position-Content Interactions in DINOv2 Block Index 6.



B. Linear Probe Evaluation for Factorization

To verify that the factorization cleanly separates positional and content structure, we train token-level linear probes to decode the spatial coordinate of each patch token from (i) the unfactorized block activations and (ii) the content factor μ_c .

Data preparation. For each model and each layer, we extract either the block activations $A_\ell(\mathbf{x}) \in \mathbb{R}^{T \times D}$ or the content factor $\mu_c(\mathbf{x}) \in \mathbb{R}^{T \times D}$ from 5,000 ImageNet images, where T is the number of patch tokens and D is the embedding dimension. We treat every patch token as an individual training example. With 5,000 images and 256 patch tokens per image, this yields $N = 5000 \times 256 = 1.28M$ token examples per probe. Each token has a discrete position label $\in \{1, \dots, 256\}$ corresponding to its grid location. We collect the token representations into a matrix \mathbf{H} and their position labels into a vector \mathbf{y} :

$$\begin{aligned}\mathbf{H} &\in \mathbb{R}^{N \times D}, \\ \mathbf{y} &\in \{1, \dots, 256\}^N, \\ N &= 5000 \times 256 = 1,280,000.\end{aligned}$$

Probe architecture and optimization. A multinomial logistic regression classifier is trained using cross-entropy loss. We use the Adam optimizer (learning rate 10^{-2} , batch size 8,192, 20 epochs), implemented in a GPU-optimized batched training loop. No nonlinearities, regularization, or data augmentation are used. After training, the probe is evaluated on all tokens. Accuracy reflects the fraction of tokens for which the predicted coordinate matches the ground-truth position. High accuracy when decoding from the unfactorized block activations confirms that spatial topology remains linearly accessible in the raw representation, whereas chance-level accuracy from the content factor μ_c verifies that positional information is successfully removed by the factorization.

C. Interaction Maps and Mode Specialization

Factor-projected mode energies (per image). For head h with SVD $\mathbf{W} = \mathbf{U}\Sigma\mathbf{V}^\top$ and mode i (columns $\mathbf{u}_i, \mathbf{v}_i$; singular value σ_i), we compute the factor-projected query/key codes ($\mu^Q, \mu^K \in \{\mu_L, \mu_P, \mu_C\}$) for image \mathbf{x} :

$$\begin{aligned}\mathbf{z}_{\cdot,i}^Q(\mathbf{x}) &= \mu^Q(\mathbf{x}) \mathbf{u}_{h,i} \in \mathbb{R}^d, \\ \mathbf{z}_{\cdot,i}^K(\mathbf{x}) &= \mu^K(\mathbf{x}) \mathbf{v}_{h,i} \in \mathbb{R}^d\end{aligned}$$

The factor-projected energy for mode i in image \mathbf{x} is then:

$$\xi_{\cdot,i}^{(i)}(\mathbf{x}) = \|\mathbf{z}_{\cdot,i}^Q(\mathbf{x}) \sigma_i \mathbf{z}_{\cdot,i}^K(\mathbf{x})\|_2^2.$$

Mode Normalization. After computing the per-image mode energies $\xi_{\cdot,i}^{(i)}(\mathbf{x})$, we normalize across all modes within the same head and interaction type, and then take the expectation across images:

$$\bar{\xi}_{\cdot,i}^{(i)} = \mathbb{E}_{\mathbf{x}} \left(\frac{\xi_{\cdot,i}^{(i)}(\mathbf{x})}{\sum_j \xi_{\cdot,j}^{(j)}(\mathbf{x})} \right), \quad \sum_i \bar{\xi}_{\cdot,i}^{(i)} = 1.$$

From nine to six interactions (symmetrization). We symmetrize directional pairs:

$$\begin{aligned}\bar{\mathcal{E}}_{P-C}^{(i)} &= \frac{1}{2} \left(\bar{\mathcal{E}}_{P-C}^{(i)} + \bar{\mathcal{E}}_{C-P}^{(i)} \right), \\ \bar{\mathcal{E}}_{L-C}^{(i)} &= \frac{1}{2} \left(\bar{\mathcal{E}}_{L-C}^{(i)} + \bar{\mathcal{E}}_{C-L}^{(i)} \right), \\ \bar{\mathcal{E}}_{L-P}^{(i)} &= \frac{1}{2} \left(\bar{\mathcal{E}}_{L-P}^{(i)} + \bar{\mathcal{E}}_{P-L}^{(i)} \right),\end{aligned}$$

yielding six undirected interactions $\{L-L, P-P, C-C, L-P, L-C, P-C\}$.

Interaction Maps For each transformer layer ℓ and each of the six symmetrized interaction types, we compute their normalized mode energies $\bar{\mathcal{E}}_{L-L}^{(i)}, \bar{\mathcal{E}}_{P-P}^{(i)}, \bar{\mathcal{E}}_{C-C}^{(i)}, \bar{\mathcal{E}}_{L-P}^{(i)}, \bar{\mathcal{E}}_{L-C}^{(i)}, \bar{\mathcal{E}}_{P-C}^{(i)}$. These energies are arranged into head-by-mode matrices, with heads (12 per layer) along the rows and singular modes (64 per head) along the columns (Appendix Figure 17 and Figure 6). Brighter values indicate modes that contribute more strongly to a given interaction after within-interaction normalization.

These maps reveal two consistent trends. First, fine-grained specialization emerges at the level of individual modes rather than at the level of whole heads: even within a single head, different modes selectively support different informational interactions. Second, the distribution of energy across modes differs systematically between ViT and DINOv2. In the supervised ViT, interaction energy is concentrated into a small subset of top singular modes, producing visibly sharper matrices. In contrast, DINOv2 exhibits a more distributed pattern, with content-driven and position-driven interactions expressed across a broader range of modes. This provides a layer-wise view of the richer mode spectrum and higher stable rank highlighted in the main text.

Ternary plots. To visualize mode specialization, we group the six undirected interaction energies for each mode i into three informational families: a *layer* family (all interactions involving the layer factor), a *position* family (interactions involving the position factor), and a *content* family (interactions involving the content factor). Mixed interactions (such as position-content or layer-content) contribute to both relevant families. We then normalize the three family scores to obtain barycentric coordinates for each mode. Each singular

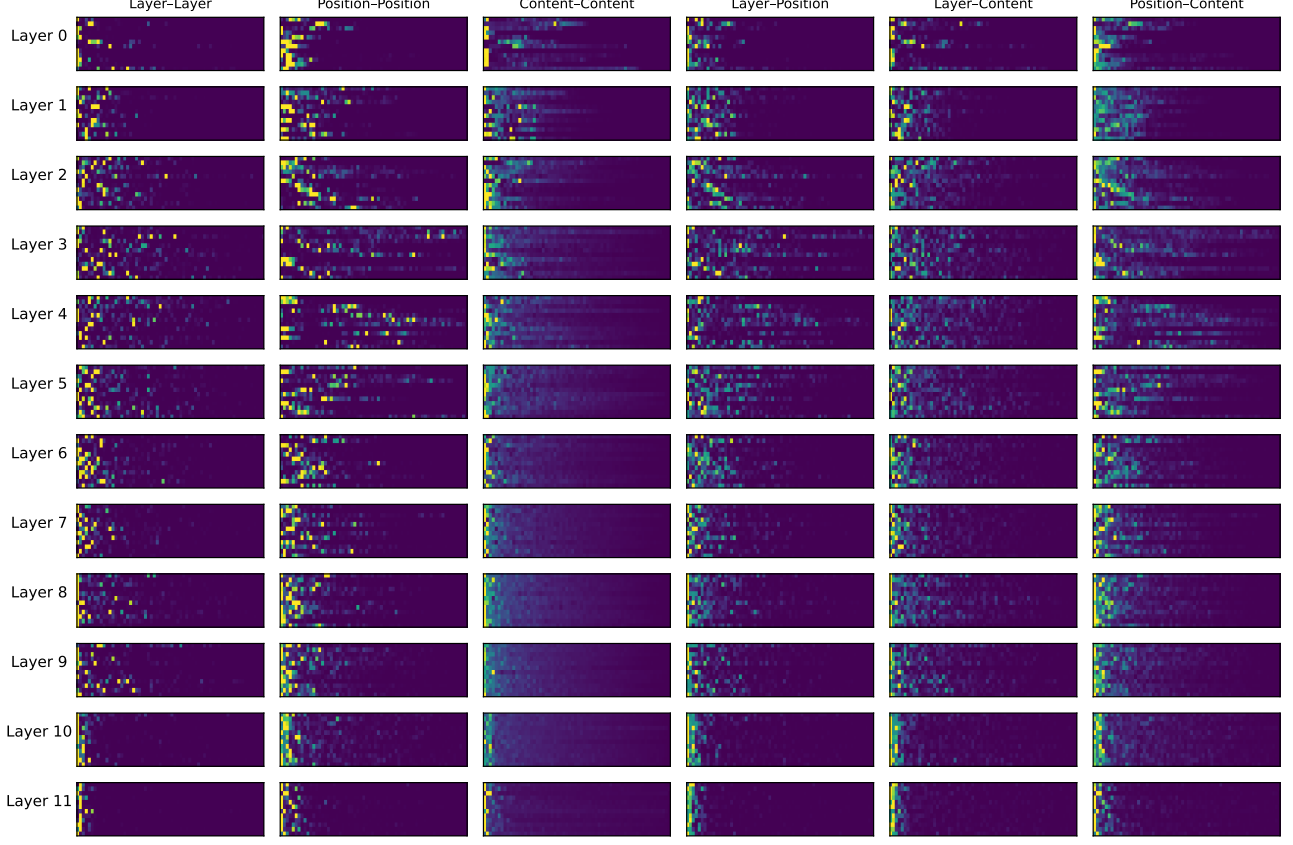


Figure 17. **Normalized Energy per Mode and Head for ViT.** Each subplot visualizes the normalized energy for a specific interaction factor (columns) at each layer (rows). For each layer and interaction, the x-axis represents the bi-orthogonal modes, sorted by singular value from highest (left) to lowest (right) and the y axis represents the head. The color intensity shows the relative contribution of a single mode to a given interaction’s total energy (normalized horizontally for each head). The visualization reveals that the energy is concentrated in the select few top singular modes.

mode i is plotted at $(\hat{S}_L^{(i)}, \hat{S}_P^{(i)}, \hat{S}_C^{(i)})$ in the ternary simplex. Points near the content vertex are dominated by content-only (resp. position-only, layer-only) interactions. Points along the content–position edge indicate genuine co-activation of these families (localization-aware semantic integration).

Layer aggregation and Density Plots. For each transformer layer ℓ , we collect all singular modes across its $H = 12$ attention heads and $K = 64$ modes per head, yielding a total of $H \times K = 768$ mode coordinates per layer:

$$\mathcal{S}_\ell = \{(\hat{S}_L^{(i,h,\ell)}, \hat{S}_P^{(i,h,\ell)}, \hat{S}_C^{(i,h,\ell)}) : i \in [1, K], h \in [1, H]\}.$$

Each point corresponds to one singular mode’s relative energy distribution among layer-, position-, and content-dominated interactions. The per-layer ternary plots (Appendix Figure 19) visualize these 768 points per layer, showing how the informational specialization of modes evolves with depth. To summarize the architecture as a

whole, we pool all modes across layers, heads, and factors ($12 \times 12 \times 64 = 9216$ points) and plot their distribution in a single ternary simplex (Figure 18), revealing each model’s overall specialization bias.

To visualize specialization trends more clearly, we convert each per-layer set of mode coordinates \mathcal{S}_ℓ into a continuous density map within the ternary simplex. Each mode’s position $(\hat{S}_L^{(i,h,\ell)}, \hat{S}_P^{(i,h,\ell)}, \hat{S}_C^{(i,h,\ell)})$ is treated as a sample from the layer’s specialization distribution. We estimate a smooth kernel density over the simplex using hexagonal binning (analogous to a 2D KDE) and normalize the resulting occupancy so that densities sum to one within each layer. This produces the color-coded plots in Figure 20, where *lighter regions* indicate a higher concentration of modes with similar informational composition. These per-layer densities make it possible to track how factor specialization sharpens or redistributes with depth—e.g., early layers show broader mixtures, while deeper layers in DINOv2 cluster tightly toward content- or position-dominated vertices.

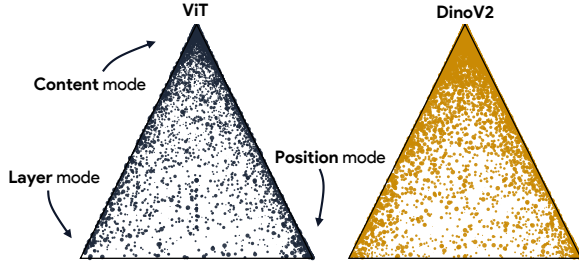


Figure 18. **Mode specialization across the network.** Each triangle (ViT on left and DINOv2 on right) visualizes all singular modes from all attention heads and all layers, plotted in barycentric coordinates indicating the relative contribution of layer, position, and content interaction families. Points near a vertex indicate modes dominated by layer-only, position-only, or content-only interactions while points along edges reflect mixed operator families (e.g., content–position). These layer-aggregated ternary plots summarize the overall computational footprint of different interactions in each model and highlight the characteristic content-centric organization induced by self-supervised training.

Spectra. We additionally visualize the singular value spectra for every head and layer (Appendix Figure 21). While the ViT spectra decay rapidly, again indicating that only a few modes dominate each head, DINOv2 retains a flatter decay profile across the entire depth of the network. This architectural difference implies that self-supervised training supports more active communication channels throughout the model. Together, these visualizations expand the analysis of [our characterization of information flow](#) by showing how informational factors are distributed within each head and how this structure evolves across depth, providing a complete head-and-mode resolution of information flow in both models.

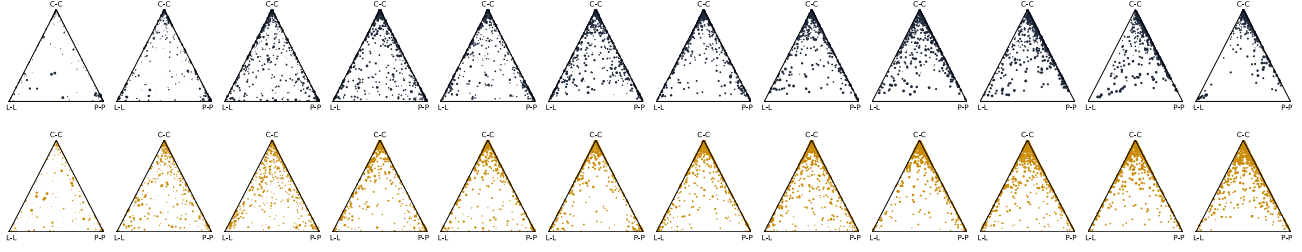


Figure 19. **Per-layer mode specialization.** Each triangle visualizes all singular modes from all attention heads in a given layer, plotted in barycentric coordinates indicating the relative contribution of layer, position, and content interaction families. Points near a vertex indicate modes dominated by layer-only, position-only, or content-only interactions while points along edges reflect mixed operator families (e.g., content–position). The top row shows ViT; the bottom row shows DINOv2. Across both models, modes exhibit strong functional specialization, with DINOv2 showing a broader distribution of content-dominated modes and more content–position operators in mid-layers.

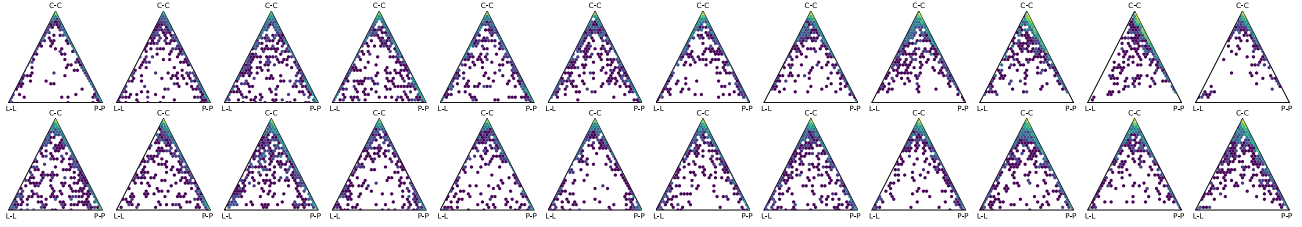


Figure 20. **Per-layer distribution of modes.** Each triangle aggregates singular modes from all heads within a layer and displays their barycentric coordinates as a hexagonal density map. Lighter regions indicate higher density. ViT (top) and DINOv2 (bottom) show progressively increasing content bias with DINOv2 showing a richer mixture of content–position operators. These density plots reveal how the computational roles of modes evolve across depth at the level of the entire layer.

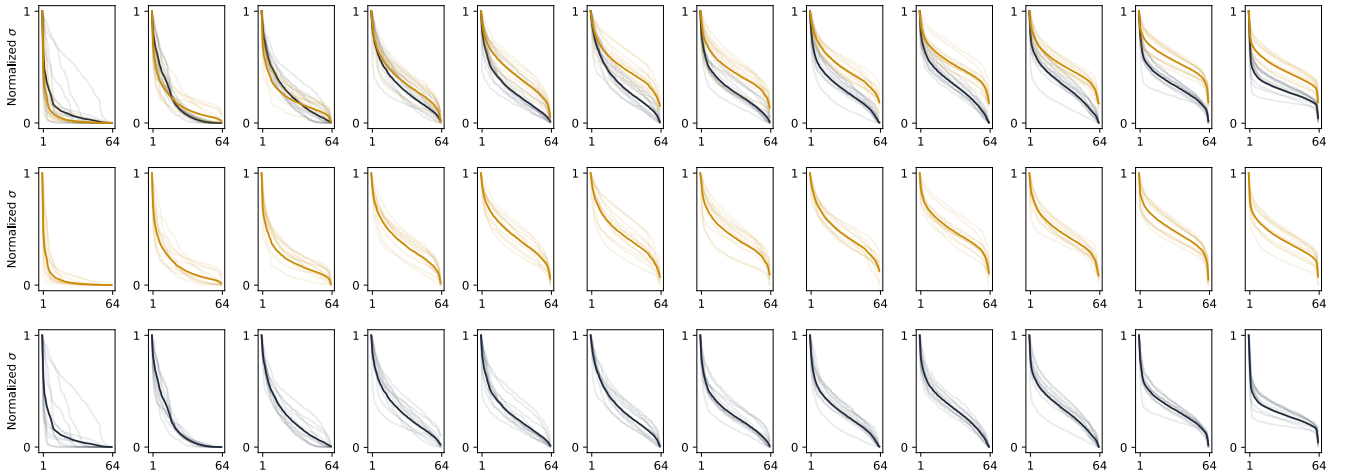


Figure 21. Details of W singular values (normalized) across layers; ViT top, DINO bottom.

D. Position representation

Figure 22 and Figure 23 provide an expanded version of the positional PCA analysis introduced in Figure 9. For each layer ℓ and each model, we estimate the positional factor $\mu_p^{(p)}$ for all spatial tokens $\in \mathbb{R}^{196 \times 768}$ for ViT-B/16 (a 14×14 grid of tokens) and $\in \mathbb{R}^{256 \times 768}$ for DINOv2-B/14 (a 16×16 grid of tokens). We then run PCA directly on $\mu_p^{(p)}$ and project each token onto the top three principal components, yielding 3D coordinates in $\mathbb{R}^{196 \times 3}$ and $\mathbb{R}^{256 \times 3}$. The resulting 3D point clouds are rendered from five viewing angles $\theta \in \{0^\circ, 45^\circ, 90^\circ, 135^\circ, 180^\circ\}$, so that the geometry of the positional manifold can be inspected in 3D.

Each point is colored solely by its 2D spatial location using a fixed position colormap (shown in the inset of Figure 9): tokens are arranged on a 14×14 or 16×16 grid and mapped to the corresponding cell in the color chart (e.g., the bottom-right token, index 256, uses the dark blue color highlighted in the legend). This color coding is identical across layers and models, allowing us to visually track how the 2D positional grid is embedded and distorted in the PCA space.

For the supervised ViT (Figure 22), the positional manifold quickly collapses with depth: by mid-layers the tokens cluster along an effectively one-dimensional curve, and the 2D grid structure becomes indistinguishable from most viewpoints. This matches the quasi-1D structure already visible in the top panel of Figure 9, and shows that the collapse is robust to viewing angle rather than an artifact of a particular projection. In contrast, DINOv2 (Figure 23) maintains an approximately planar 2D sheet across all layers: rotating the view reveals that neighboring tokens remain embedded in a smooth, grid-like surface, and no comparable 1D collapse occurs even at depth. These complementary visualizations therefore substantiate our claim that self-supervised training preserves an explicit 2D positional scaffold, whereas supervised training compresses positional geometry to a much lower effective dimensionality.

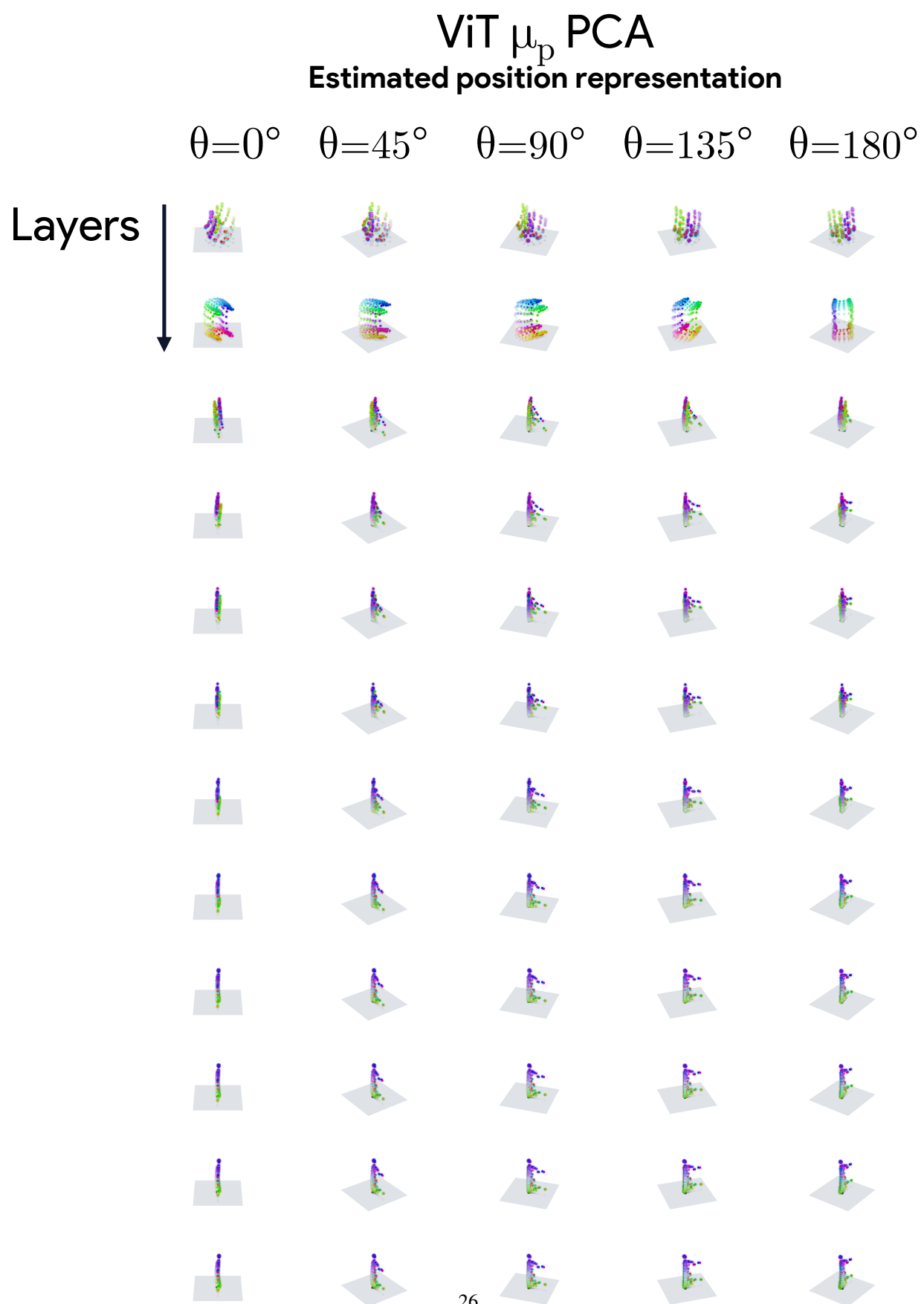


Figure 22. Estimated Position (μ_p) in ViT visualized via 3D PCA under 5 rotation angles from layer 1 (top) to layer 12 (bottom).

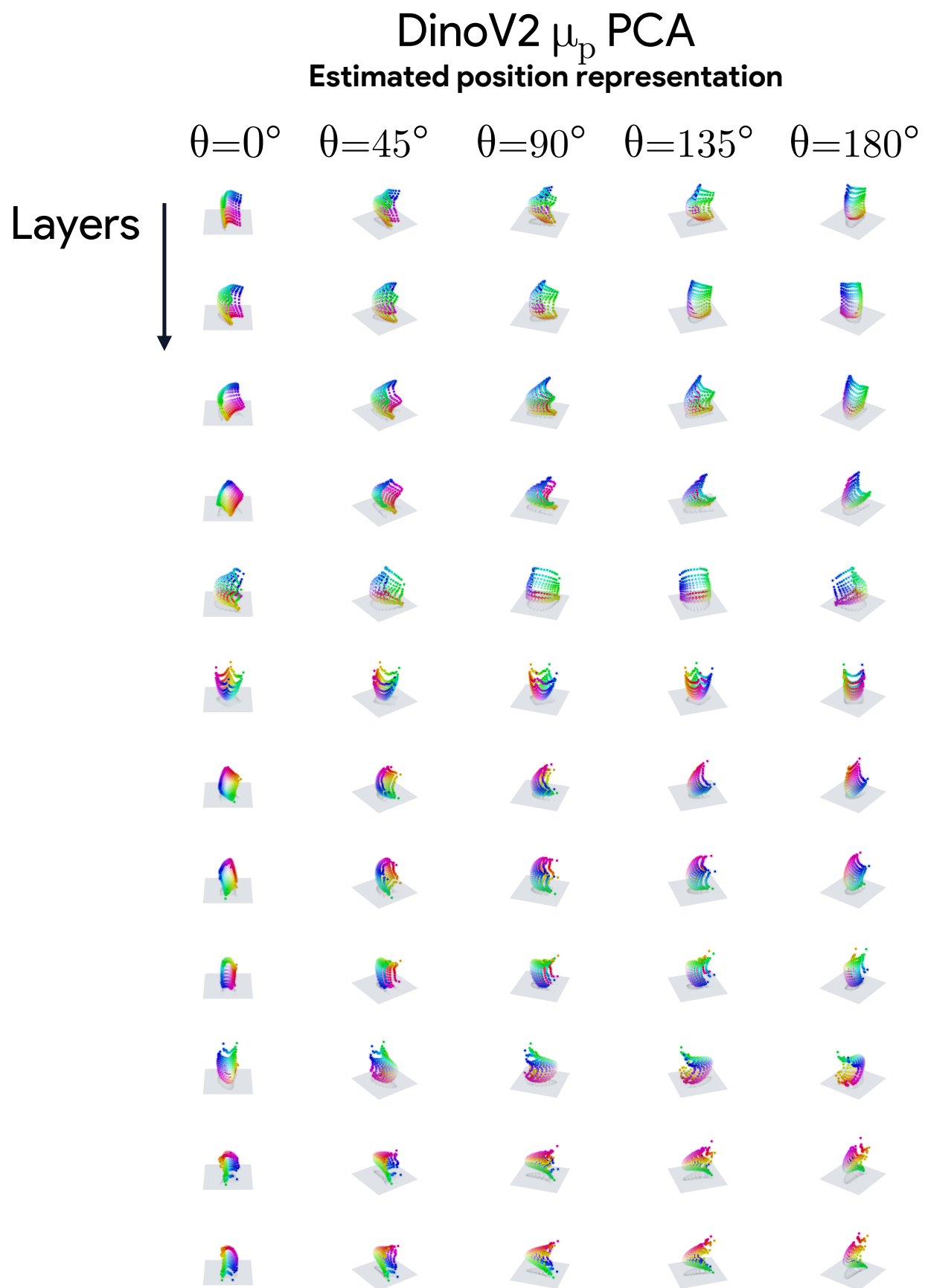


Figure 23. Estimated Position (μ_p) in DINOv2 visualized via 3D PCA under 5 rotation angles from layer 1 (top) to layer 12 (bottom).

E. Enrichment of Content Across Layers

To quantify how token representations evolve across depth, we compute pairwise layer–layer Pearson correlations between the flattened representations, both for the unfactorized block activations and for the content factor. This analysis characterizes how the geometry of patch representations changes across layers.

For each model, layer $\ell \in \{0, \dots, L-1\}$, and image \mathbf{x} , we extract either the unfactorized block activations $A_\ell(\mathbf{x}) \in \mathbb{R}^{T \times D}$ or the content factor $\mu_{c,\ell}(\mathbf{x}) \in \mathbb{R}^{T \times D}$, where T is the number of total tokens (patch and the special tokens) and D is the embedding dimension. We then flatten the representation into a vector $\mathbf{h}_\ell(\mathbf{x}) \in \mathbb{R}^{TD}$ and compute all pairwise Pearson correlations between these vectors, yielding an $L \times L$ matrix. Averaging these matrices over all $N = 5,000$ images produces a single layer–layer similarity matrix for each model. High correlations along the diagonal indicate representational stability across adjacent layers, whereas high off-diagonal correlations indicate that distant layers preserve similar token content. Lower off-diagonal values instead reflect progressive transformation of the representations with depth. Empirically, DINOv2 exhibits smoother mid-layer evolution than supervised ViT, consistent with the richer content operations revealed by our mode-level analysis.