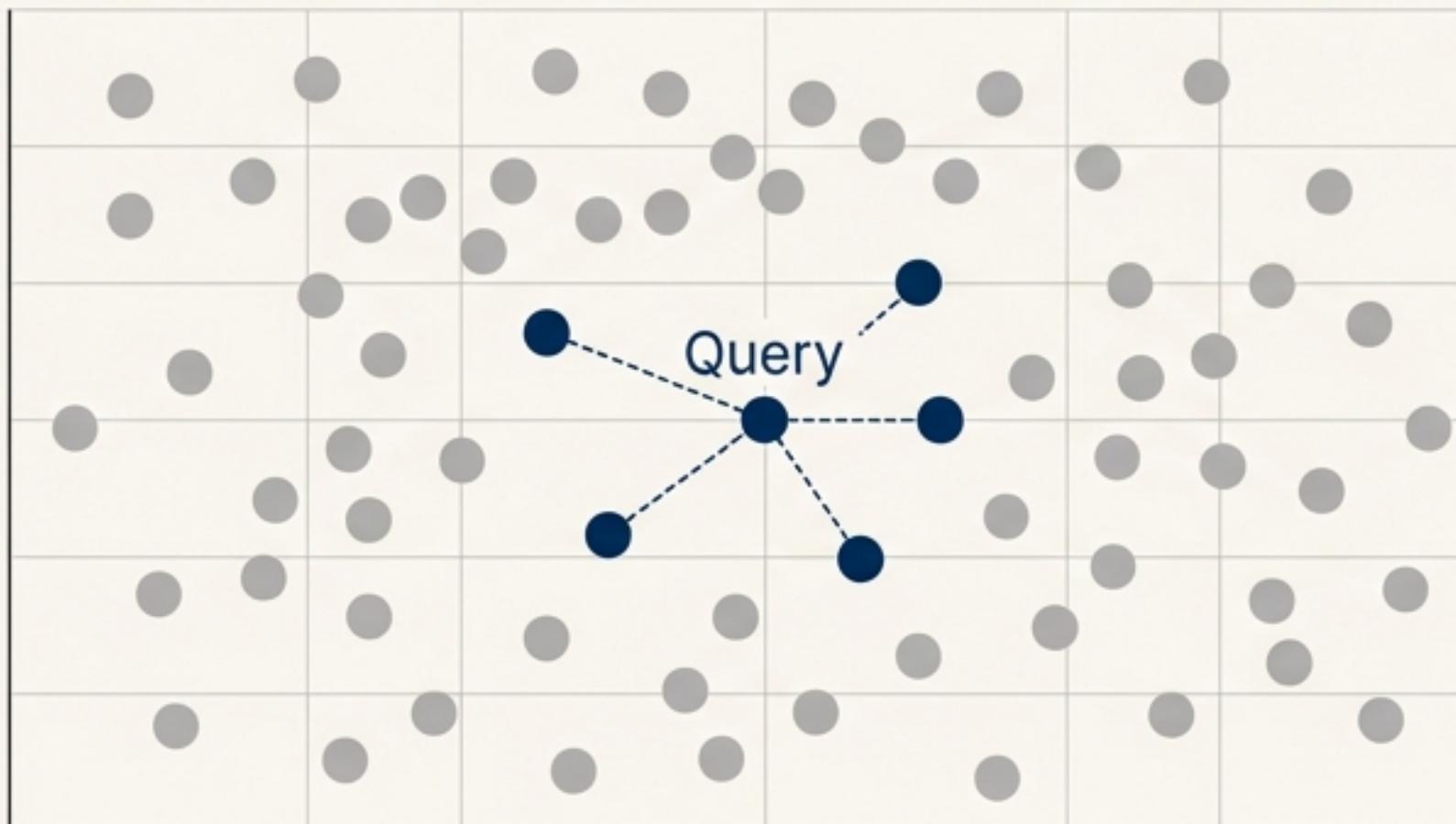


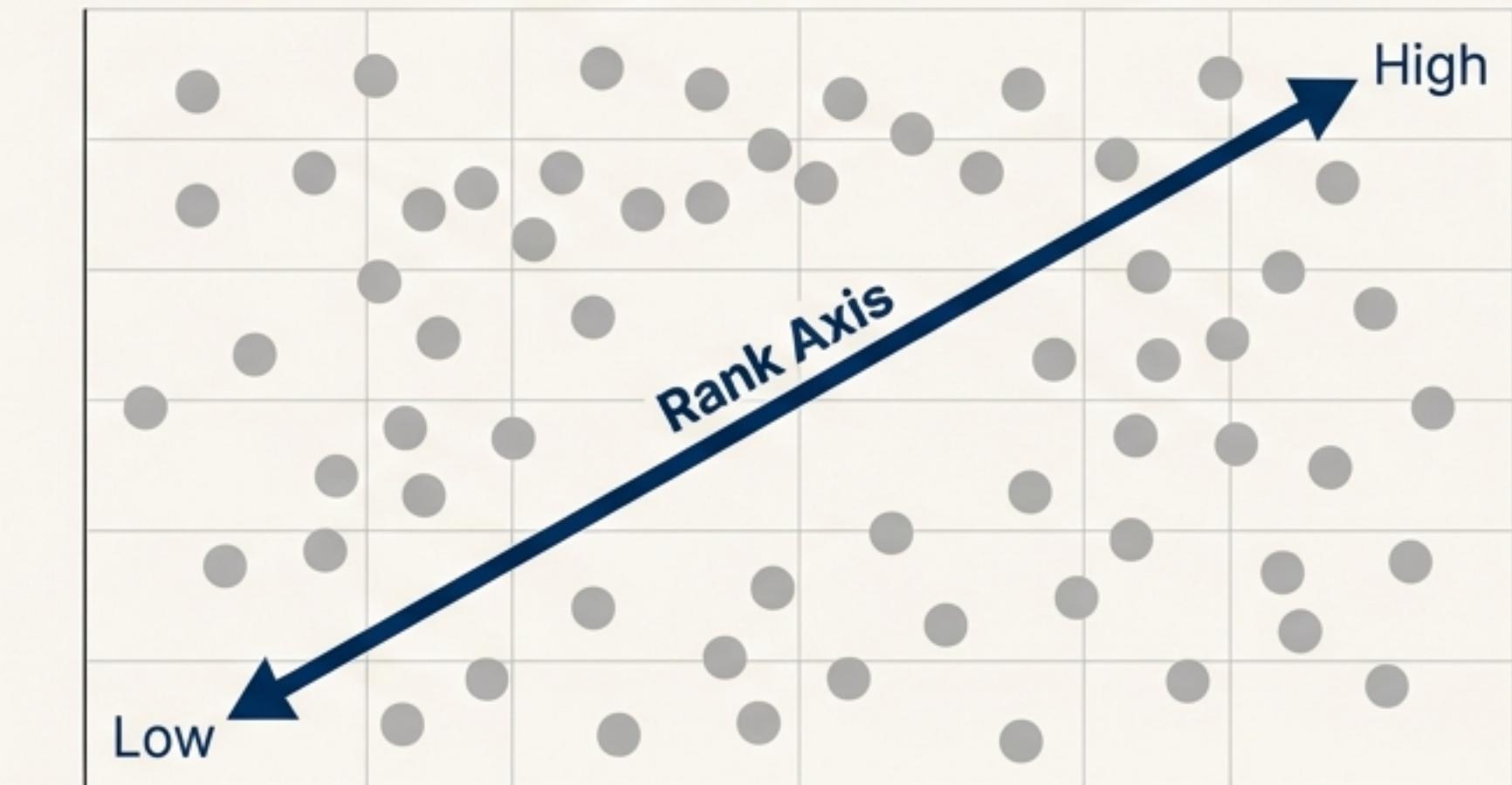
# We Thought Embeddings Were Just for Similarity. A Recent Paper Claims They Are Also Ordered.

- **The Conventional View:** Vector databases use embeddings for *local similarity search*. Given a query image, we find its nearest neighbors in the embedding space. This is the foundation of modern image retrieval.
- **A Provocative Hypothesis:** This paper (“On the rankability of visual embeddings”) argues for a different, hidden property: *global ordering*. They claim that embeddings possess an inherent, linear structure that allows for ranking along continuous attributes.
- **Our Investigation Today:** As scientists, we must be skeptical. Is this claim credible? How robust is it? And is it even practical? Let's put it to the test.

**Retrieval: Local Similarity**



**Ranking: Global Ordering**



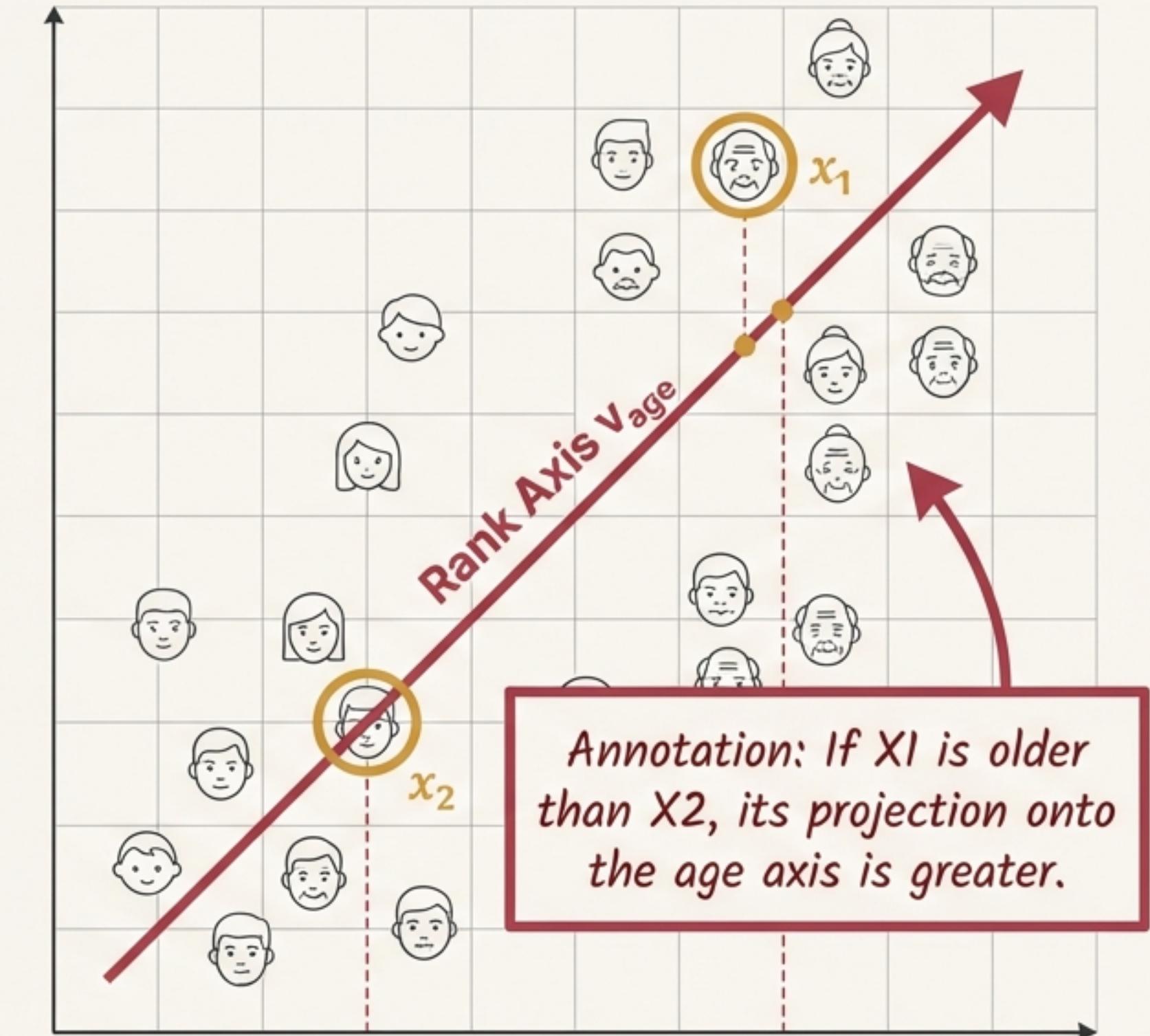
# What Exactly is “Rankability”?

The paper defines rankability as follows:

- An embedding model  $f$  is **rankable** for a continuous attribute  $A$  (e.g., age) if there exists a **rank axis**  $v_A$  (a direction vector in the embedding space).
- ...such that for any two images,  $x_1$  and  $x_2$ , if  $A(x_1) \geq A(x_2)$ , then the projection of their embeddings onto the axis preserves this order:  $v_A^T f(x_1) \geq v_A^T f(x_2)$ .

## In Simple Terms:

If we can find a single direction in the embedding space that sorts all images by a specific attribute (like youngest to oldest), the embedding is rankable for that attribute.



# The Experimental Setup: How Do You Measure Something Like This?

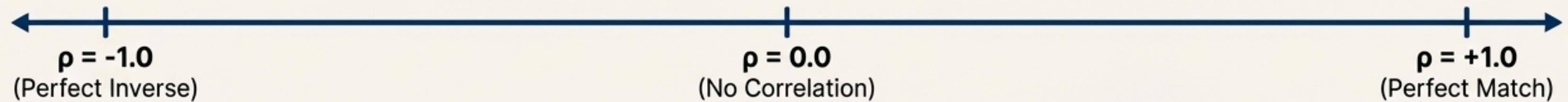
The authors' approach is straightforward but requires careful benchmarking.

## 1. Finding the Rank Axis

They train a simple **linear regressor** on a set of image embeddings and their corresponding attribute labels (e.g., age). The weights of this regressor define the rank axis  $v_A$ .

## 2. Measuring Performance

They use **Spearman's Rank Correlation Coefficient** (SRCC /  $\rho$ ). This metric measures how well the predicted ranking (from projecting onto  $v_A$ ) matches the true attribute ranking.



## 3. Contextualizing the Score (The Baselines)

A raw SRCC score is meaningless without context. They compare against three crucial reference points:



### Lower Bound (No-train)

SRCC from an untrained, randomly initialized encoder. The "no information" baseline.

### Upper Bound (Nonlinear)

SRCC from a 2-layer MLP regressor. Estimates total ordinal information available in the embedding.

### Upper Bound (Finetuning)

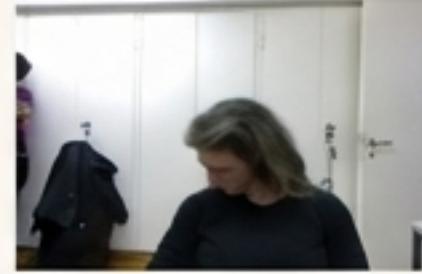
SRCC after finetuning the entire encoder. Shows the architectural capacity for the attribute.

# A First Look at the Evidence: Visualizing the Rank Axes

Let's look at what happens when images are sorted along the rank axes found by the linear regressor (using CLIP-ViT-B/32 embeddings). The numbers below each image are the ground-truth attributes.

## Head Pose

(looking down → looking up)



0th percentile, -55



25th percentile, -10



50th percentile, 14



75th percentile, 36



100th percentile, 74

## Crowd Count

(not crowded → crowded)



0th percentile, 129



25th percentile, 233



50th percentile, 751



75th percentile, 944



100th percentile, 1408

## Aesthetics

(bad → good)



0th percentile, 25



25th percentile, 48



50th percentile, 58



75th percentile, 78



100th percentile, 76

## Age

(young → old)



0th percentile, 5



25th percentile, 18



50th percentile, 35



75th percentile, 55



100th percentile, 80

# The Quantitative Verdict: Visual Embeddings Are Generally Rankable

## Key Findings

The visual results are convincing, but let's examine the numbers. This table shows the SRCC scores, averaged across 7 different encoder architectures.

- **Key Finding:** The 'Rankability' scores (from the linear regressor) are not just high (mostly 0.6 to 0.9), they are consistently **much closer to the upper bounds** (Nonlinear, Finetuned) than to the "no-train" lower bound.
- **Example (Age on Adience):** The rankability score of 0.861 is only slightly below the nonlinear upper bound of 0.878, but far above the lower bound of 0.266.
- **Interpretation:** This implies that most of the ordinal information present in the embeddings is already structured *linearly*. A complex nonlinear model isn't needed to extract it.

## SRCC ( $\rho$ ) Averaged Across 7 Architectures

Attribute (Dataset)	No-train lower bound	Rankability (main)	Nonlinear upper bound	Finetuned upper bound
Age (UTKFace)	0.199	0.766	0.776	0.799
Age (Adience)	0.266	0.861	0.878	0.910
Crowd (UCF-QNRF)	0.220	0.843	0.854	0.886
Aesthetics (AVA)	0.341	0.696	0.716	0.773
Crowd (JHU-CROWD)	0.203	0.782	0.801	0.844
Aesthetics (KonIQ)	0.352	0.727	0.748	0.791
Time (Flickr2018)	0.089	0.645	0.677	0.712

# But Not All Embeddings Are Created Equal

Digging deeper, we find that the rankability depends heavily on the encoder architecture.

- **Main Observation:** CLIP-based embeddings are generally more rankable than their non-CLIP counterparts across most attributes.
- **Interesting Exception:** DINOv2 massively outperforms CLIP on head yaw, suggesting some properties are better captured by self-supervised methods.
- **Conclusion:** While rankability is a general property, the multimodal, language-aligned pre-training of CLIP appears to embed ordinal concepts more linearly and robustly.

## Rankability (SRCC) by Encoder Architecture

	Non-CLIP Models				CLIP Models		
	RN50	ViTB32	CNX	DINO-B14	CLIP-RN50	CLIP-ViT32	CLIP-CNX
<b>Age</b> (UTKFace)	0.698	0.714	0.741	0.739	0.764	0.779	0.783
<b>Age</b> (Adience)	0.839	0.846	0.871	0.859	0.899	0.907	0.928
Crowd (UCF-QNRF)	0.806	0.818	0.823	0.835	0.849	0.861	0.869
<b>Aesthetics</b> (AVA)	0.613	0.631	0.644	0.660	0.722	0.741	0.750
Crowd (JHU-CROWD)	0.737	0.748	0.761	0.772	0.791	0.805	0.815
<b>Aesthetics</b> (KonIQ)	0.689	0.701	0.710	0.719	0.737	0.749	0.758
<b>Time</b> (Flickr2018)	0.564	0.583	0.599	0.614	0.655	0.677	0.696
<b>Recency</b> (HCI)	0.598	0.615	0.631	0.650	0.789	0.807	0.820
<b>Yaw</b> (Kinect)	0.321	0.349	0.380	0.804	0.398	0.423	0.440

Exception!

# And Not All Attributes Are Easily Ranked

Even with the best models, some attributes are clearly harder to rank than others.

- **High Rankability:** Attributes like Age, Crowd Count, and Head Pitch are ranked very well across most models ( $\text{SRCC} > 0.8$ ).
- **Poor Rankability:** Head Yaw and especially Head Roll show very poor linear rankability ( $\text{SRCC}$  of 0.434 and 0.218 on average).
- **Hypothesis (from the authors):** The strength of an attribute's rank axis is likely proportional to the variety and presence of that attribute in the model's original, large-scale pre-training data.
- **Takeaway:** The linear structure exists, but its signal strength varies. It is not a universal constant for all possible concepts.

**Rankability (SRCC) by Encoder Architecture**

		Non-CLIP Models				CLIP Models		
		RN50	ViTB32	CNX	DINO-B14	CLIP-RN50	CLIP-ViTB32	CLIP-CNX
High Rankability	Pitch (Kinect)	0.810	0.830	0.838	0.833	0.846	0.868	0.874
	Age (Adience)	0.634	0.706	0.740	0.737	0.768	0.878	0.760
	Crowd (UCF-QNRF)	0.700	0.763	0.761	0.877	0.768	0.784	0.745
Lower Rankability	Aesthetics (AVA)	0.321	0.349	0.380	0.804	0.398	0.423	0.440
	Crowd (JHU-CROWD)	0.400	0.344	0.353	0.802	0.401	0.487	0.480
	Aesthetics (KonIQ)	0.304	0.431	0.450	0.418	0.484	0.494	0.498
Poor Rankability	Roll (Kinect)	0.185	0.210	0.230	0.220	0.201	0.225	0.245
	Time (Flickr2018)	0.408	0.413	0.418	0.483	0.468	0.512	0.504
	Recency (HCI)	0.217	0.204	0.220	0.249	0.201	0.225	0.245

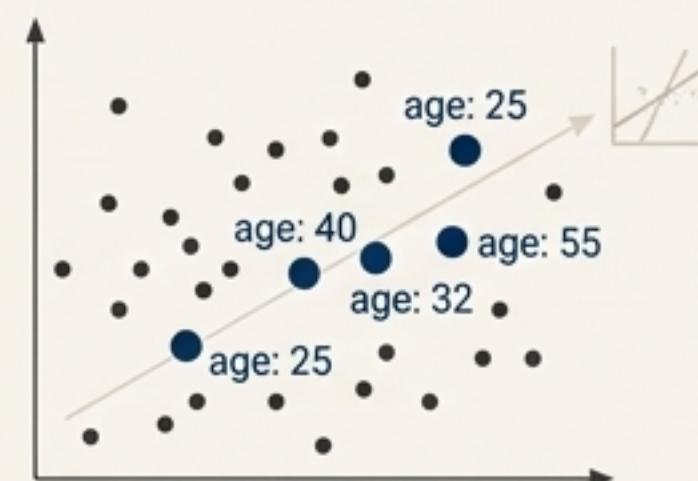
# This is Interesting... But Finding the Axis Requires Full Supervision. Can We Do Better?

So far, we've found the rank axis by training a regressor on a full dataset with expensive, continuous attribute labels. This is not practical for real-world applications where a user wants to rank by an arbitrary new attribute. The paper investigates **two more** efficient, low-supervision alternatives:

## 1. Few-Shot Learning with Labels

The standard approach.

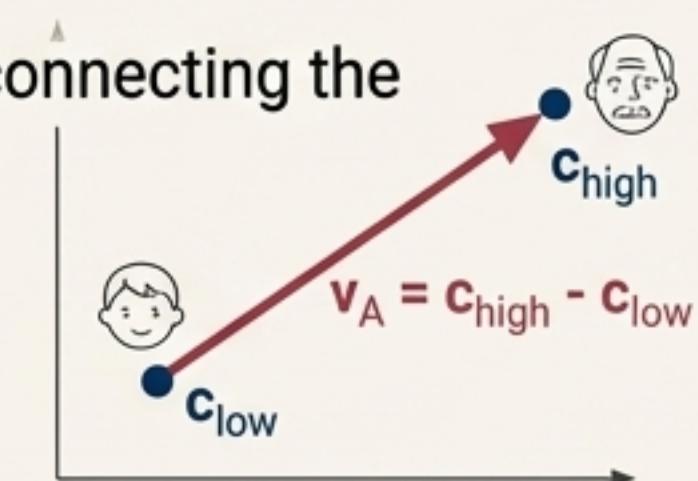
Train the linear regressor, but with only a small fraction (e.g., 2, 4, 8... samples) of the labeled data.



## 2. Steering Vector with Extreme Pairs (No Labels)

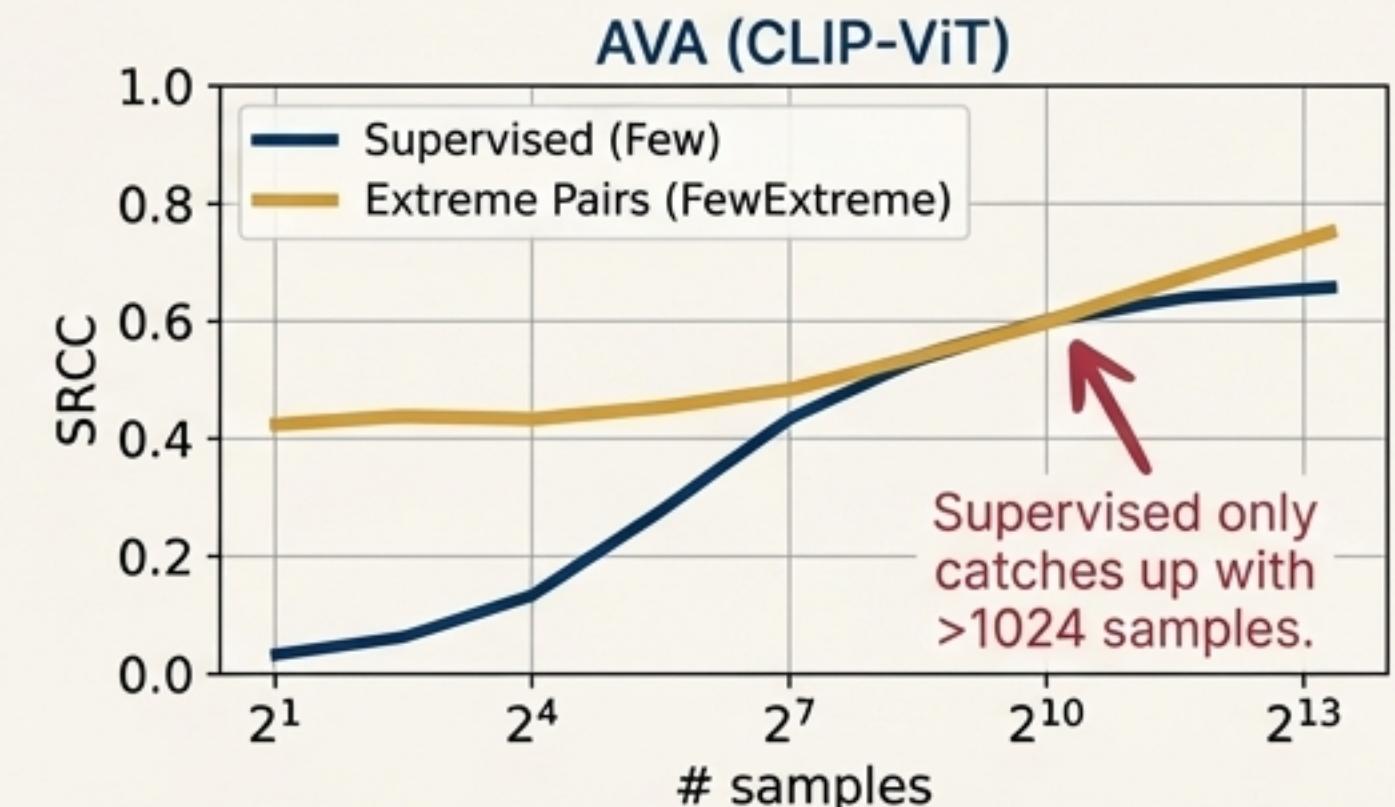
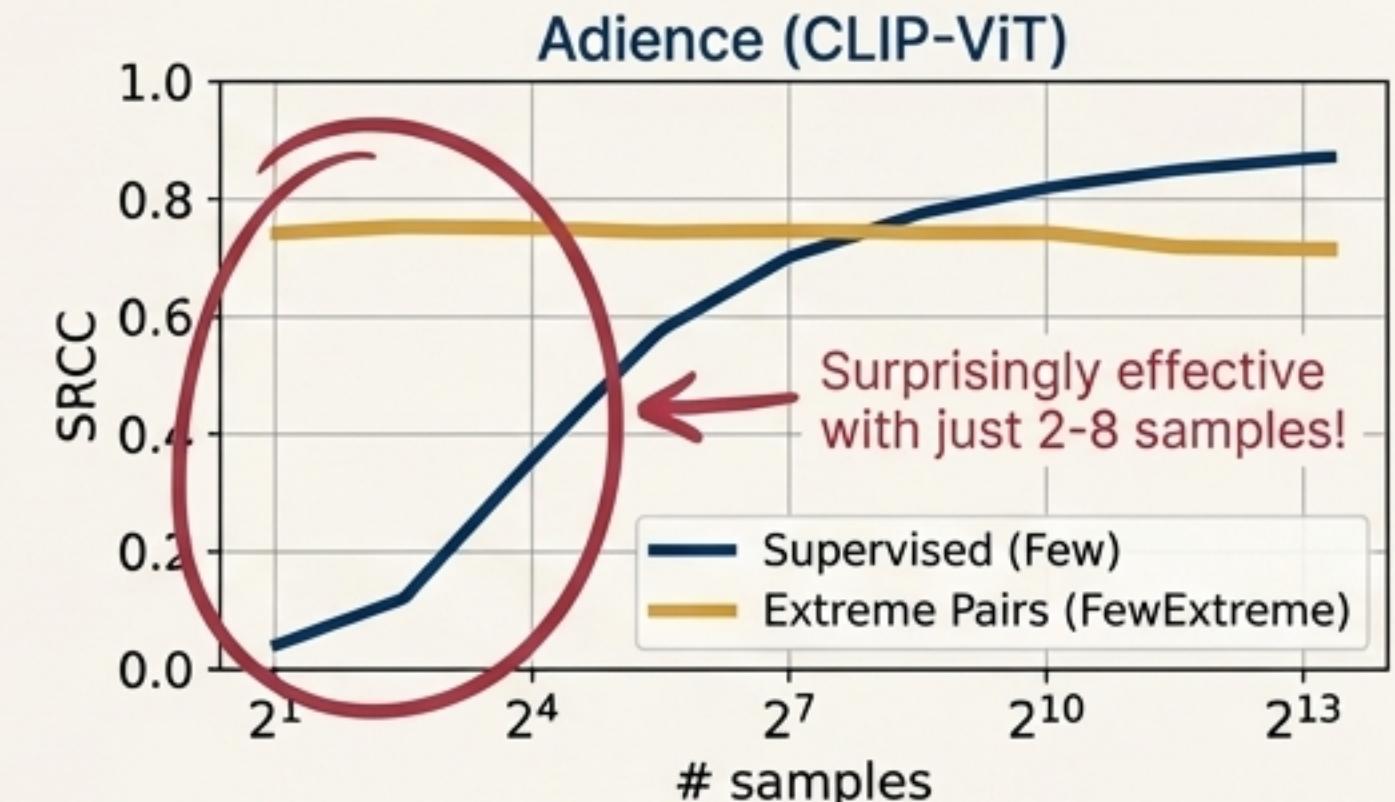
A more clever, learning-free approach.

- Find a few images from the extreme ends of the desired attribute (e.g., one image of a child, one of an elderly person).
- Compute the centroid embedding for the 'low' examples ( $\mathbf{c}_{\text{low}}$ ) and 'high' examples ( $\mathbf{c}_{\text{high}}$ ).
- The rank axis is simply the vector connecting the centroids:  $\mathbf{v}_A = \mathbf{c}_{\text{high}} - \mathbf{c}_{\text{low}}$ .



# In the Low-Data Regime, Extreme Pairs Beat Supervised Learning

- The results of comparing the two efficient methods are striking.
- **Key Observation:** When the number of available samples is very small (up to ~128), the Extreme Pairs method consistently outperforms few-shot learning with full labels.
- On the Adience dataset, a rank axis from just *two* extreme images (one young, one old) achieves an SRCC of ~0.75. The supervised method with two labeled samples is near zero.
- **Practical Takeaway:** For quick, on-the-fly ranking, finding two good “extreme” examples is a more effective use of resources than meticulously labeling a few dozen random samples.



# How Universal is a Rank Axis? Testing Transferability

If concepts like “age” have an inherent linear structure, an axis learned on one dataset should work on another. The authors tested this by training on one dataset and evaluating on another for the same attribute.

## Findings:

### 1. Transfer is Non-Trivial and Asymmetric:

- An “age” axis trained on Adience transfers to UTKFace with a respectable SRCC of 0.68.
- However, the reverse is worse: UTKFace -> Adience only yields an SRCC of 0.55.

### 2. Rank Directions are Correlated:

- The cosine similarity between the “age” axes from UTKFace and Adience is 0.36. While far from 1.0, this is a significant alignment.

## Conclusion:

There appears to be a “universal” direction for core attributes like age, but it is warped by dataset-specific biases. The axes are related, but not identical.

**Visual 1: SRCC Transfer Matrix (Age)**

		Evaluated on	
		UTKFace	Adience
Trained on	UTKFace	0.81	0.55
	Adience	0.68	0.91

**Visual 2: SRCC Transfer Matrix (Crowd Count)**

		Evaluated on	
		UCF-QNRF	ST-A
Trained on	UCF-QNRF	0.82	
	ST-A	0.73	

**Visual 3: Cosine Similarity**

Cosine Similarity( $v_{age\_UTKFace}$ ,  $v_{age\_Adience}$ ) = **0.36**

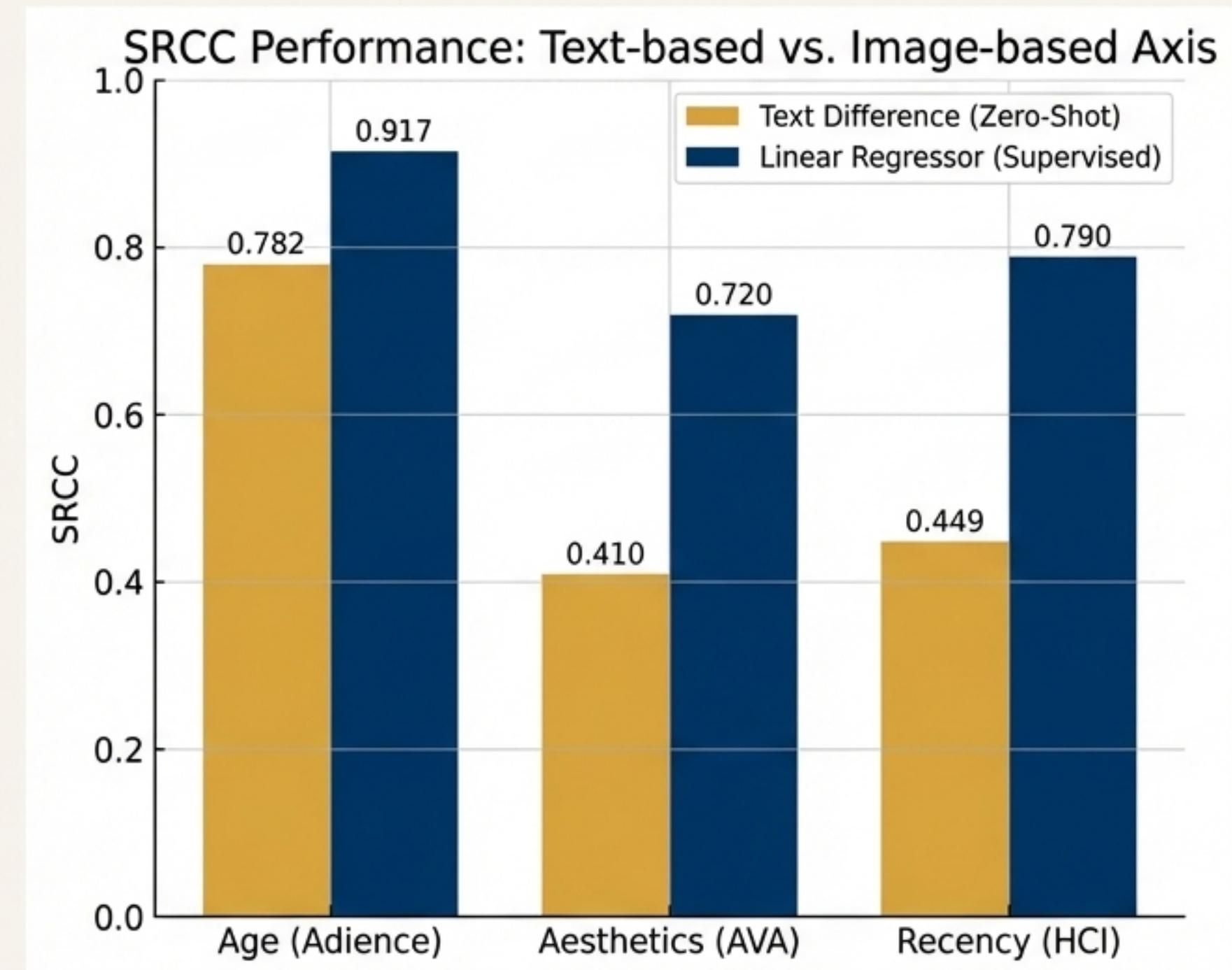
# The Final Frontier: Can We Find the Axis with Zero-Shot Text Prompts?

For models like CLIP, could we find the rank axis without any images at all, using just text? The authors tested two text-based methods:

1. **Single Prompt**: The axis is the embedding of a single prompt, e.g., "a photo of an old person".
2. **Text Difference**: The axis is the embedding difference between two prompt embeddings, e.g., "a photo of an old person" – "a photo of a young person".

**The Verdict:** Language-based prompting lags considerably behind image-based methods.

**Takeaway:** While language provides a directional hint, the true geometric rank axis is best defined by visual examples. The modality gap between text and image embeddings remains a significant factor.



# What This Investigation Uncovered: The Broader Implications

Our deep dive confirms the paper's central claim and reveals some profound implications for how we understand and use embeddings.



## 1. Embeddings Have an Interpretable Structure

The existence of linear, ordinal axes suggests the latent space is not an indecipherable black box. We can "probe" it along meaningful, human-understandable directions. This is a step towards explainability.



## 2. New Possibilities for Vector Databases

This property could simplify database operations. Instead of complex, dedicated hybrid indexing structures for filtering on metadata, we might be able to perform ranking *directly* on the vectors themselves, as long as the attribute is encoded in a rank axis.

### Example Use Case

A photo app lets a user sort selfies by "age appearance." It computes the `v\_age` axis using just two reference images (a child, an elderly person) and ranks the entire album without any extra metadata.

# But Let's Be Clear About the Limitations

This “rankability” is a powerful property, but it is not a magic bullet. We must remember its boundaries.



**It Applies to Continuous, Ordinal Attributes:** The method works for concepts that have a natural order (e.g., age, crowd size, aesthetics).



**It Does Not Handle Categorical Attributes:** There is no “reasonable ordering” between discrete categories like “cat” and “dog.” An embedding space can’t have a single linear axis that sorts animals by species.



**Hybrid Indexing Is Not Obsolete:** For filtering on categorical or complex metadata, dedicated database structures are still necessary. Rankability supplements, but does not replace, traditional attribute filtering.



**Performance Varies:** As we saw, the strength of the rank axis varies significantly by model and by attribute. It must be validated for each specific use case.

# A Final Takeaway: Embeddings Possess a Surprisingly Linear, Rankable Structure

The core discovery of this work is that modern visual embeddings, especially those from models like CLIP, don't just cluster similar images—they organize them along linear, ordinal axes. This is an emergent property that is both unexpected and practically useful. We've gone from seeing embeddings as a "bag of points" to a structured, geometric space.



## The Next Great Question

“Could one potentially also characterise embeddings as interpretable collections of latent ordinal subspaces?” – The Authors

This is an exciting and open direction for future work—deconstructing the hidden geometry of these powerful models.