

# Contents

<b>1. Introduction</b>	<b>1</b>
<b>2. Preliminaries</b>	<b>2</b>
2.1. Interpretable-by-design classification models	2
2.2. Sparse coding for concept extraction . . . . .	2
2.3. Geometric Structures of Meanings in Vector Embeddings . . . . .	3
<b>3. Hierarchical Concept Embedding Model</b>	<b>3</b>
3.1. Well-clustered synset embeddings . . . . .	3
3.2. Hierarchical Orthogonality . . . . .	4
<b>4. Hierarchical Concept Pursuit</b>	<b>4</b>
4.1. Hierarchical Dictionary Construction . . . . .	4
4.2. Hierarchical Orthogonal Matching Pursuit . . . . .	5
<b>5. Experiments</b>	<b>6</b>
5.1. Synthetic Experiments . . . . .	6
5.2. Real-World Experiments . . . . .	7
<b>6. Conclusion</b>	<b>8</b>
<b>A Extended Related Work</b>	<b>12</b>
<b>B Preliminaries</b>	<b>12</b>
B.1. Interpretable-by-design vs. Chain-of-Thought	12
B.2. Canonical regular simplex . . . . .	13
<b>C Limitations</b>	<b>13</b>
<b>D Proofs</b>	<b>13</b>
D.1. Proof of Prop. 3.1 . . . . .	13
D.2. Proof of Prop. 3.2 . . . . .	14
D.3. Proof of Prop. 3.3 . . . . .	14
D.4. Intermediate results for Prop. 4.1 . . . . .	14
D.5. Proof of Prop. 4.1 . . . . .	15
<b>E Step-by-step construction of a Hierarchical Concept Embedding</b>	<b>15</b>
E.1. Feasible Subspace Induced by Hierarchical Orthogonality . . . . .	16
<b>F. Additional Experimental Results</b>	<b>16</b>
F.1. Additional Synthetic Experiment Details . . . . .	16
F.2. Text Interpretation of Synset Differences . . . . .	17
F.3. Additional Real-Data Experiment Details . . . . .	17

## A. Extended Related Work

**Interpretable-by-design models.** Interpretable-by-design models aim to provide explanations for their predictions by

using human-interpretable concepts as intermediate representations. Early works explored attribute-based classification for face verification [34] and learning to detect unseen object classes through attribute transfer [36]. Subsequent works include Concept Activation Vectors [28], which use linear classifiers to identify directions in the embedding space corresponding to specific concepts. Concept Bottleneck Models [29] extend this idea by training a model to predict concepts before predicting the final output. In an adjacent line of work, Information Pursuit [20] is used as a criterion to choose the most relevant concepts [7, 9, 30]. More recent works have explored leveraging pre-trained embeddings and sparse coding for identifying specific concept directions [4, 8, 19]. Our work builds upon these foundations by introducing a concept embedding framework that captures the hierarchical relationships among synsets in interpretable image classification.

**Sparse Recovery.** Sparse recovery aims to recover a sparse signal from a set of observations, often using techniques such as Orthogonal Matching Pursuit (OMP) [54] and Basis Pursuit [11]. Sparse coding has been widely used in image processing [43, 44], signal processing, and machine learning. Although there have been works on hierarchical sparse coding [25, 27, 35], they do not consider the hierarchical structure of concepts in the context of interpretable models or deep representation learning. More recently, sparse autoencoders (SAEs) [6, 13, 48] have been used in vision-language models [12, 50, 70] to interpret the structure of concepts beyond sparsity. Unlike SAE-based approaches, our work focuses on structured sparse coding with an explicitly hierarchical dictionary derived from semantic synset relationships.

## B. Preliminaries

### B.1. Interpretable-by-design vs. Chain-of-Thought

Chain-of-Thought (CoT) prompting [67] is a technique that enables large language models to generate step-by-step reasoning traces before producing a final answer, often improving performance on complex reasoning tasks. However, using CoT as an interpretability method does not guarantee faithfulness<sup>5</sup> because the final prediction is conditioned on both the input and the generated chain of thoughts. Several recent audits further question the faithfulness of CoT explanations [3, 38, 65]. Therefore, the development of CoT does not make interpretable-by-design models obsolete. In fact, there is a trend toward making foundation models (e.g., LLMs or diffusion models) interpretable-by-design by enforcing an interpretable concept bottleneck in the latent space [23, 60].

<sup>5</sup>An explanation is *faithful* if it accurately reflects the true computation process to the final prediction [24].

## B.2. Canonical regular simplex

Define

$$\tilde{\mathbf{s}}_j = \mathbf{e}_j - \frac{1}{b} \mathbf{1}, \quad j = 1, \dots, b, \quad (9)$$

where  $\{\mathbf{e}_j\}_{j=1}^b$  are the standard basis vectors of  $\mathbb{R}^b$  and  $\mathbf{1} \in \mathbb{R}^b$  is the all-ones vector. These centred vertices satisfy

$$\sum_{j=1}^b \tilde{\mathbf{s}}_j = \mathbf{0} \text{ and } \tilde{\mathbf{s}}_j^\top \tilde{\mathbf{s}}_k = \begin{cases} 1, & j = k, \\ -\frac{1}{b-1}, & j \neq k, \end{cases} \text{ i.e. they}$$

form a regular  $(b-1)$ -simplex of unit edge length in  $\mathbb{R}^{b-1}$ .

## C. Limitations

While our framework demonstrates clear advantages in concept recovery, it also has several limitations:

**Embedding Constraints.** Our theoretical analysis (Prop. 3.3) establishes that embedding a hierarchy with leaf depth  $L$  and branching ratio  $b$  requires ambient dimension  $d \geq L + b - 1$ . For deep hierarchies (large  $L$ ) or highly branching structures (large  $b$ ), this constraint becomes restrictive. Real-world embeddings from models such as CLIP typically have fixed dimensions (e.g.,  $d = 768$ ), which limits the depth and complexity of hierarchies that can be faithfully represented. Moreover, as hierarchies deepen, the half-angles of the cones containing each subtree (Prop. 3.2) must decrease geometrically. As mentioned in §3.1, this limitation may necessitate exploring alternative geometries, such as hyperbolic spaces [15, 49], for more faithful hierarchical representations. At the same time, by using pretrained embeddings, HCEP provides accurate and interpretable classification without the additional compute required to finetune large models, and pretrained models are available and improving across multiple domains, making HCEP easy to extend. That said, hierarchy-aware finetuning could further improve the geometric conditions in §3 and thus boost HCEP’s performance; we leave this direction to future work.

**Hierarchy Quality Dependence.** The performance of Hierarchical OMP critically depends on the quality of the pre-defined hierarchy. For ImageNet-based datasets, we leverage the well-curated WordNet hierarchy, which provides semantically meaningful relationships. However, for CIFAR-100, we rely on taxonomy induction methods [71], which may produce hierarchies with inconsistencies or unclear relationships. That said, high-quality hierarchies already exist in many domains beyond common objects, e.g., RadLex [37] for radiology, DERM12345 [69] for skin lesions, and iNaturalist [66] for species classification. Moreover, the quality of LLM-based taxonomy induction is steadily improving, as demonstrated by our CIFAR-100 experiments (§5.2). When the hierarchy is noisy, HCEP may produce

incorrect explanation paths. A promising direction for future work is to interpolate between hierarchical and non-hierarchical solutions via  $\ell_1$  minimization with a hierarchy regularizer [25], thereby allowing the method to degrade gracefully when hierarchy quality is uncertain.

**Computational Complexity.** Hierarchical OMP with beam search (Alg. 2) has complexity  $O(TBb|\mathcal{D}_{\text{active}}|)$ , where  $T$  is the number of iterations,  $B$  is the beam width,  $b$  is the branching ratio, and  $|\mathcal{D}_{\text{active}}|$  is the size of the active dictionary at each level. In contrast, OMP has complexity  $O(T|\mathcal{D}|)$ , where  $|\mathcal{D}|$  is the size of the entire dictionary. For large branching ratios or deep hierarchies, this may become computationally expensive. While we demonstrate better concept recovery accuracy over standard OMP, the computational cost remains a practical consideration for deployment at scale. In practice, the beam search hypotheses can be parallelized on GPU, which significantly reduces wall-clock time overhead (see Fig. 10).

**Outlier Handling.** HCEP assumes that the input belongs to a leaf class in the provided hierarchy. For outliers within a known class (e.g., a cat with three legs), classical results on the robustness of sparse coding to corruption [68] suggest that the outlier will still be closest to the correct synset path, while a larger sparse reconstruction error can flag such cases. If a novel class is entirely absent from the hierarchy, taxonomy induction can be used to append the new class before applying HCEP.

## D. Proofs

### D.1. Proof of Prop. 3.1

**Statement.** If subtree containment (Eq. (3)) and sibling-cone disjointness (Eq. (4)) hold, then the subtrees rooted at sibling nodes do not overlap.

*Proof.* Suppose  $j$  and  $j'$  are distinct children of a parent node  $i$ . We show that their subtrees are disjoint.

Let  $k \in \text{desc}(j)$ . By subtree containment in Eq. (3),

$$\angle(\mathbf{a}^{(j)}, \mathbf{a}^{(k)}) \leq \theta_{\text{lev}(j)}. \quad (10)$$

Assume for contradiction that  $k \in \text{desc}(j')$  as well. Applying Eq. (3) to  $j'$  gives

$$\angle(\mathbf{a}^{(j')}, \mathbf{a}^{(k)}) \leq \theta_{\text{lev}(j')}. \quad (11)$$

Now consider the spherical triangle formed by the unit vectors  $\mathbf{a}^{(j)}/\|\mathbf{a}^{(j)}\|$ ,  $\mathbf{a}^{(k)}/\|\mathbf{a}^{(k)}\|$ , and  $\mathbf{a}^{(j')}/\|\mathbf{a}^{(j')}\|$ . The triangle inequality yields

$$\begin{aligned} \angle(\mathbf{a}^{(j)}, \mathbf{a}^{(j')}) &\leq \angle(\mathbf{a}^{(j)}, \mathbf{a}^{(k)}) + \angle(\mathbf{a}^{(k)}, \mathbf{a}^{(j')}) \\ &\leq \theta_{\text{lev}(j)} + \theta_{\text{lev}(j')}, \end{aligned} \quad (12)$$

which contradicts sibling-cone disjointness in Eq. (4). Therefore  $k$  cannot belong to both  $\text{desc}(j)$  and  $\text{desc}(j')$ .

Since  $j$  and  $j'$  were arbitrary siblings, the subtrees rooted at sibling nodes are disjoint.  $\square$

## D.2. Proof of Prop. 3.2

**Statement.** If the half-angles satisfy  $\theta_{l+1} \leq \min\{r, 1/b\}\theta_l$  with  $r \in (0, 1/2)$ , then there exists a placement of the nodes such that *subtree containment* (Eq. (3)) and *sibling-cone disjointness* (Eq. (4)) hold.

*Proof.* We verify the two requirements separately.

**Step 1: Subtree containment.** A sufficient condition for Eq. (3) is that the cumulative half-angles of all lower levels do not exceed the half-angle budget at level  $l$ , namely

$$\sum_{k=l+1}^L \theta_k \leq \theta_l, \quad \forall i \in \{1, \dots, N_L\}, \quad l = \text{lev}(i). \quad (13)$$

Assume first that the half-angles satisfy the geometric decrease

$$\theta_{l+1} \leq r \theta_l. \quad (14)$$

Then, for any level  $l$ ,

$$\sum_{k=l+1}^L \theta_k \leq \sum_{k=1}^{L-l} \theta_l r^k \text{ (by Eq. (14))} \quad (15)$$

$$= \theta_l \sum_{k=1}^{L-l} r^k \quad (16)$$

$$= \theta_l r \sum_{k=0}^{L-l-1} r^k \quad (17)$$

$$\leq \theta_l r \frac{1 - r^{L-l-1}}{1 - r} \text{ (sum of a geometric series).} \quad (18)$$

Taking  $L \rightarrow \infty$  yields

$$\theta_l r \frac{1 - r^{L-l-1}}{1 - r} \rightarrow \theta_l \frac{r}{1 - r} \leq \theta_l, \quad r < \frac{1}{2},$$

so Eq. (13) holds. This proves subtree containment.

**Step 2: Sibling-cone disjointness.** Next, we show that Eq. (4) follows from the simpler sufficient condition

$$\theta_{l+1} \leq \frac{1}{b} \theta_l. \quad (19)$$

Consider any 2D plane containing the parent cone axis. In that plane, a cone of half-angle  $\alpha$  appears as a planar angle of magnitude  $2\alpha$ . Therefore, one can place  $b$  child cone axes inside the parent cone without overlap whenever  $b$  child angles of size  $2\theta_{l+1}$  fit inside the parent angle  $2\theta_l$ , i.e., whenever  $b\theta_{l+1} \leq \theta_l$ . This is exactly Eq. (19).

Since the assumption of the proposition is  $\theta_{l+1} \leq \min\{r, 1/b\}\theta_l$ , both Step 1 and Step 2 hold simultaneously. Hence there exists a placement satisfying subtree containment and sibling-cone disjointness.  $\square$

## D.3. Proof of Prop. 3.3

**Statement.** Under the hierarchical orthogonality constraints in Eq. (30) and the regular-simplex difference condition in Eq. (31) at every internal node up to depth  $L - 1$  of a hierarchy whose leaves are at depth  $L$  in ambient space  $\mathbb{R}^d$ , it is necessary that the ambient dimension satisfies the depth–dimension condition  $d \geq L + b - 1$ .

*Proof.* Fix a level  $l \geq 0$  and consider the ancestor path  $\mathbf{A}_l = \{\mathbf{a}^{(\pi_0)}, \dots, \mathbf{a}^{(\pi_l)}\}$ . The hierarchical orthogonality constraints in Eq. (30) define the affine feasible set for child candidates:

$$\mathcal{V}_l = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{A}_l^\top \mathbf{x} = \mathbf{h}_l\},$$

which is Eq. (34). When the ancestor vectors are linearly independent, as is generically the case because each level introduces a new non-collinear direction, we have

$$\dim \mathcal{V}_l = d - (l + 1).$$

Now apply the regular-simplex condition in Eq. (31). It requires placing  $b$  child points whose differences relative to a feasible origin in  $\mathcal{V}_l$  form a regular  $(b-1)$ -simplex. Since such a simplex has affine hull of dimension  $b-1$ , it can be embedded in  $\mathcal{V}_l$  only if

$$\dim \mathcal{V}_l \geq b - 1.$$

Combining the two displays yields, for every level  $l$ ,  $d - (l + 1) \geq b - 1 \iff d \geq l + b$ . Requiring this inequality to hold at the deepest internal level  $l = L - 1$  gives  $d \geq (L - 1) + b = L + b - 1$ , as claimed.  $\square$

## D.4. Intermediate results for Prop. 4.1

**Lemma D.1** (Column normalization equivalence). *Let  $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_k] \in \mathbb{R}^{d \times k}$  with arbitrary nonzero column norms ( $\|\mathbf{d}_j\|_2 > 0$  for all  $j \in [k]$ ), and define the diagonal matrix  $\mathbf{W} := \text{diag}(\|\mathbf{d}_1\|_2, \dots, \|\mathbf{d}_k\|_2)$  and the column-normalized dictionary  $\widehat{\mathbf{D}} := \mathbf{D}\mathbf{W}^{-1}$ . For any  $s$ -sparse  $\mathbf{z} \in \mathbb{R}^k$  with support  $S$ , set  $\widehat{\mathbf{z}} := \mathbf{W}\mathbf{z}$ . Then  $\mathbf{x} = \mathbf{D}\mathbf{z} = \widehat{\mathbf{D}}\widehat{\mathbf{z}}$  and  $\text{supp}(\widehat{\mathbf{z}}) = S$ . Moreover, OMP run on  $\mathbf{D}$  with the selection rule*

$$j^* \in \arg \max_j \frac{|\langle \mathbf{r}, \mathbf{d}_j \rangle|}{\|\mathbf{d}_j\|_2} = \arg \max_j \frac{|\langle \mathbf{r}, \mathbf{d}_j \rangle|}{\|\mathbf{d}_j\|_2 \|\mathbf{r}\|_2} \quad (20)$$

*is identical (same index picked at every iteration) to OMP run on  $\widehat{\mathbf{D}}$  with the usual (unnormalized) correlation rule. Equivalently, this selects the atom with the highest absolute cosine similarity to the residual.*

*Proof.* The representation identity is immediate:

$$\mathbf{x} = \mathbf{D}\mathbf{z} = \mathbf{D}\mathbf{W}^{-1}\mathbf{W}\mathbf{z} = \widehat{\mathbf{D}}\widehat{\mathbf{z}}, \quad (21)$$

and clearly  $\text{supp}(\hat{\mathbf{z}}) = S$  because  $\mathbf{W}$  is diagonal with positive entries.

For the selection rule, note that  $\hat{\mathbf{d}}_j = \mathbf{d}_j / \|\mathbf{d}_j\|_2$ , so

$$\langle \mathbf{r}, \hat{\mathbf{d}}_j \rangle = \frac{\langle \mathbf{r}, \mathbf{d}_j \rangle}{\|\mathbf{d}_j\|_2}. \quad (22)$$

Thus maximizing correlation with  $\hat{\mathbf{d}}_j$  is exactly the same as maximizing normalized correlation, equivalently absolute cosine similarity, with  $\mathbf{d}_j$ .

Finally, diagonal rescaling of the columns in  $\mathbf{D}_S$  does not change  $\text{span}(\mathbf{D}_S)$ , so the orthogonal projector onto this span is invariant under the rescaling  $\mathbf{D}_S \mapsto \hat{\mathbf{D}}_S$ . Therefore, once the same index is selected, both versions of OMP produce the same least-squares fit and hence the same residual at every step. By induction over the iterations, the selected indices are identical throughout the run.  $\square$

**Definition D.2** (Exact Recovery Coefficient (ERC) on normalized dictionary).<sup>6</sup> For a support  $S$  with  $\hat{\mathbf{D}}_S$  full column rank, define

$$\text{ERC}(\hat{\mathbf{D}}; S) := \max_{j \in S^c} \|\hat{\mathbf{D}}_S^\dagger \hat{\mathbf{d}}_j\|_1, \quad (23)$$

$$\text{ERC}(\hat{\mathbf{D}}; S | T) := \max_{j \in T \setminus S} \|\hat{\mathbf{D}}_S^\dagger \hat{\mathbf{d}}_j\|_1, \quad (24)$$

for any  $T \supseteq S$ , where  $\|\cdot\|_1$  denotes the vector  $\ell_1$  norm. Here  $S^c := [k] \setminus S$ , so  $T \setminus S = T \cap S^c \subseteq S^c$ ; thus  $\text{ERC}(\hat{\mathbf{D}}; S | T)$  is the same maximum as  $\text{ERC}(\hat{\mathbf{D}}; S)$ , but over the smaller index set  $T \setminus S$ .

**Lemma D.3** (Monotone ERC improvement under subtree restriction). *Let  $\mathbf{D} \in \mathbb{R}^{d \times k}$  have arbitrary nonzero column norms and let  $\mathbf{z}$  be  $s$ -sparse with support  $S$ . Let  $T_0 \supset T_1 \supset \dots \supset T_L$  be a nested sequence with  $T_0 = [k]$  and  $S \subseteq T_\ell$  for all  $\ell = 0, \dots, L$ . Assume  $\hat{\mathbf{D}}_S$  has full column rank. Then the ERC decreases monotonically along the restriction:*

$$\text{ERC}(\hat{\mathbf{D}}; S | T_L) \leq \text{ERC}(\hat{\mathbf{D}}; S | T_{L-1}) \quad (25)$$

$$\leq \dots \leq \text{ERC}(\hat{\mathbf{D}}; S | T_0) \quad (26)$$

$$:= \text{ERC}(\hat{\mathbf{D}}; S). \quad (27)$$

*Proof.* By definition,  $\text{ERC}(\hat{\mathbf{D}}; S | T) = \max_{j \in T \setminus S} \|\hat{\mathbf{D}}_S^\dagger \hat{\mathbf{d}}_j\|_1$ . Since  $\hat{\mathbf{D}}_S^\dagger$  is fixed, shrinking  $T$  only restricts the index set over which this maximum is taken, so the value cannot increase.  $\square$

**Lemma D.4** (ERC threshold implies restricted OMP success). *Under the assumptions of Lemma D.3, if*

<sup>6</sup>We use Exact Recovery Coefficient (ERC) for the quantity whose threshold at 1 is the classical exact recovery condition in Tropp [64].

$\text{ERC}(\hat{\mathbf{D}}; S | T_L) < 1$ , then OMP run on  $\mathbf{D}$  with the normalized selection rule

$$j^* \in \arg \max_j \frac{|\langle \mathbf{r}, \mathbf{d}_j \rangle|}{\|\mathbf{d}_j\|_2} = \arg \max_j \frac{|\langle \mathbf{r}, \mathbf{d}_j \rangle|}{\|\mathbf{d}_j\|_2 \|\mathbf{r}\|_2}, \quad (28)$$

restricted to  $T_L$ , recovers  $S$  in  $s$  steps. Equivalently, OMP on  $\hat{\mathbf{D}}_{T_L}$  with the standard rule succeeds in  $s$  iterations.

*Proof.* Lemma D.1 shows that the normalized-selection rule on  $\mathbf{D}$  matches standard OMP on  $\hat{\mathbf{D}}$ . The classical noiseless ERC theorem of Tropp [64] applied to the restricted dictionary  $\hat{\mathbf{D}}_{T_L}$  then yields exact support recovery in  $s$  iterations whenever  $\max_{j \in T_L \setminus S} \|\hat{\mathbf{D}}_S^\dagger \hat{\mathbf{d}}_j\|_1 < 1$ , equivalently whenever  $\text{ERC}(\hat{\mathbf{D}}; S | T_L) < 1$ .  $\square$

## D.5. Proof of Prop. 4.1

**Statement.** There exist instances with  $\text{ERC}(\hat{\mathbf{D}}; S) \geq 1$  yet  $\text{ERC}(\hat{\mathbf{D}}; S | T_L) < 1$  for some nested  $T_0 \supset T_1 \supset \dots \supset T_L$  satisfying the right-subtree assumption  $S \subseteq T_\ell$ . Consequently, hierarchical OMP yields a strictly larger ERC-certified success region than global OMP on the full dictionary.

*Proof.* Choose any instance for which the maximizer of  $\text{ERC}(\hat{\mathbf{D}}; S)$  lies outside  $T_L$ . Removing that maximizer from the admissible index set strictly decreases the maximum, so

$$\text{ERC}(\hat{\mathbf{D}}; S | T_L) < \text{ERC}(\hat{\mathbf{D}}; S).$$

Hence it is possible to have

$$\text{ERC}(\hat{\mathbf{D}}; S) \geq 1 \quad \text{but} \quad \text{ERC}(\hat{\mathbf{D}}; S | T_L) < 1.$$

Whenever this occurs, Lemma D.4 certifies exact recovery for Hierarchical OMP on the restricted set  $T_L$ , while the standard ERC guarantee for global OMP on the full dictionary does not apply.  $\square$

## E. Step-by-step construction of a Hierarchical Concept Embedding

Assume we are at depth  $l > 0$  of the hierarchy. The path from the root to the *current parent*  $\pi_l$  (with embedding  $\mathbf{a}^{(\pi_l)} \in \mathbb{R}^d$ ) consists of the  $l+1$  ancestor vectors

$$\mathbf{A}_l = \{\mathbf{a}^{(\pi_0)}, \mathbf{a}^{(\pi_1)}, \dots, \mathbf{a}^{(\pi_l)}\}, \quad \pi_0 < \pi_1 < \dots < \pi_l. \quad (29)$$

We seek embeddings  $\{\mathbf{a}^{(j)}\}_{j=1}^b \subset \mathbb{R}^d$  for its  $b$  children that satisfy the following conditions:

(i) **Hierarchical Orthogonality:**

$$\begin{aligned} (\mathbf{a}^{(j)} - \mathbf{a}^{(\pi_k)})^\top \mathbf{a}^{(\pi_k)} &= 0, & k &= 0, \dots, l, \\ & & j &= 1, \dots, b. \end{aligned} \quad (30)$$

By induction, the current parent node  $\pi_l$  already satisfies these orthogonality constraints with respect to its ancestors.

(ii) **Regular  $(b-1)$ -simplex structure:**

$$(\mathbf{a}^{(j)} - \mathbf{g}_l)^\top (\mathbf{a}^{(k)} - \mathbf{g}_l) = \begin{cases} \lambda_l^2, & j = k, \\ -\frac{\lambda_l^2}{b-1}, & j \neq k, \end{cases} \quad (31)$$

where  $\mathbf{g}_l$  is any point satisfying all  $l+1$  equations in Eq. (30), and  $\lambda_l$  scales the simplex so that  $\angle(\mathbf{a}^{(j)}, \mathbf{a}^{(\pi_l)}) = \theta_l$ .

(iii) **Cone condition w.r.t. the current parent:**

$$\begin{aligned} \angle(\mathbf{a}^{(j)}, \mathbf{a}^{(\pi_l)}) &= \theta_l \\ \iff \|\mathbf{a}^{(j)} - \mathbf{a}^{(\pi_l)}\| &= \|\mathbf{a}^{(\pi_l)}\| \tan \theta_l, \\ j &= 1, \dots, b. \end{aligned} \quad (32)$$

This equivalence holds because Eq. (30) implies  $\langle \mathbf{a}^{(j)} - \mathbf{a}^{(\pi_l)}, \mathbf{a}^{(\pi_l)} \rangle = 0$ .

## E.1. Feasible Subspace Induced by Hierarchical Orthogonality

Let

$$\mathbf{A}_l = [\mathbf{a}^{(\pi_0)} \quad \mathbf{a}^{(\pi_1)} \quad \dots \quad \mathbf{a}^{(\pi_l)}] \in \mathbb{R}^{d \times (l+1)}. \quad (33)$$

The  $l+1$  hyperplanes in (30) intersect in the affine subspace

$$\mathcal{V}_l = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{A}_l^\top \mathbf{x} = \mathbf{h}_l\}, \quad (34)$$

where  $\mathbf{h}_l = [\|\mathbf{a}^{(\pi_0)}\|^2, \dots, \|\mathbf{a}^{(\pi_l)}\|^2]^\top$ . If the ancestor columns of  $\mathbf{A}_l$  are linearly independent<sup>7</sup>, then

$$\dim \mathcal{V}_l = d - (l+1). \quad (35)$$

To be able to embed a regular  $(b-1)$ -simplex for all depths we therefore require the *depth-dimension condition*

$$d \geq L + b. \quad (36)$$

Equation (36) quantifies the depth–dimension trade-off: one ambient degree of freedom is lost per additional ancestor constraint, while  $(b-1)$  directions are always needed to accommodate the regular simplex of conditionally independent children.

We now give a constructive procedure for the children of node  $\pi_l$  at level  $l$ .

1. **Find one feasible origin.** Solve the linear system  $\mathbf{A}_l^\top \mathbf{x} = \mathbf{h}_l$  to obtain any particular solution  $\mathbf{g}_l \in \mathcal{V}_l$ . If we also impose the cone half-angle condition, the feasible set can be further restricted to the intersection of

<sup>7</sup>This is typical because every level adds a new non-collinear vector.

$\mathcal{V}_l$  with the cone centered at the parent embedding  $\mathbf{a}^{(\pi_l)}$  and half-angle  $\theta_l$  from §3.

A convenient choice, which preserves the full half-angle budget, is to take the current parent embedding itself as the feasible origin, i.e.,  $\mathbf{g}_l = \mathbf{a}^{(\pi_l)}$ . By construction, every node is orthogonal to all of its ancestors. Hence, for each  $k \in \{0, \dots, l\}$ ,

$$(\mathbf{a}^{(\pi_l)} - \mathbf{a}^{(\pi_k)})^\top \mathbf{a}^{(\pi_k)} = 0 \quad (37)$$

$$\implies \mathbf{a}^{(\pi_l)\top} \mathbf{a}^{(\pi_k)} = \|\mathbf{a}^{(\pi_k)}\|^2. \quad (38)$$

2. **Basis for the difference linear space.** Compute an orthonormal basis

$$\mathbf{U}_l \in \mathbb{R}^{d \times (d-l-1)}, \quad (39)$$

$$\mathbf{A}_l^\top \mathbf{U}_l = \mathbf{0}, \quad (40)$$

$$\mathbf{U}_l^\top \mathbf{U}_l = \mathbf{I}_{d-l-1}, \quad (41)$$

e.g., by taking the  $d - (l+1)$  left singular vectors of  $\mathbf{A}_l$  associated with the smallest singular values, denoted  $\mathbf{U}_{l+2:d}$ .

3. **Canonical regular simplex in  $\mathbb{R}^{b-1}$ .** Use the centred construction  $\tilde{\mathbf{d}}_i = \mathbf{e}_i - \frac{1}{b}\mathbf{1}$ ,  $i = 1, \dots, b$  (cf. Eq. (9)).

4. **Scale  $\tilde{\mathbf{d}}_j$  to satisfy the cone condition.**

$$\lambda_l = \|\mathbf{a}^{(\pi_l)}\| \tan \theta_l, \quad (42)$$

$$\mathbf{d}_j = \lambda_l \tilde{\mathbf{d}}_j. \quad (43)$$

5. **Embed and translate.**

$$\mathbf{a}^{(j)} = \mathbf{a}^{(\pi_l)} + \mathbf{U}_{l+2:d} \mathbf{d}_j, \quad j = 1, \dots, b. \quad (44)$$

Choosing the scale factor

$$\lambda_l = \|\mathbf{a}^{(\pi_l)}\| \tan \theta_l, \quad (45)$$

forces  $\|\mathbf{a}^{(j)} - \mathbf{a}^{(\pi_l)}\| = \|\mathbf{a}^{(\pi_l)}\| \tan \theta_l$  for all  $j$ , and therefore  $\angle(\mathbf{a}^{(j)}, \mathbf{a}^{(\pi_l)}) = \theta_l$ , which is precisely the requirement in Eq. (32).

## F. Additional Experimental Results

### F.1. Additional Synthetic Experiment Details

Branching ratio  $b = 3$ , hierarchy depth  $L = 7$ , dimension  $d = 50$ . Initial cone half angle = 85 degrees. Initial vector norm = 0.8. Geometric reduction factor = 0.4. Total leaf nodes = 2187. Total nodes = number of atoms = 3280. Gaussian noise for each leaf for data generation  $\sigma^2 = 10^{-5}$ . Generate 5 samples per leaf for a total of 10,935 samples.

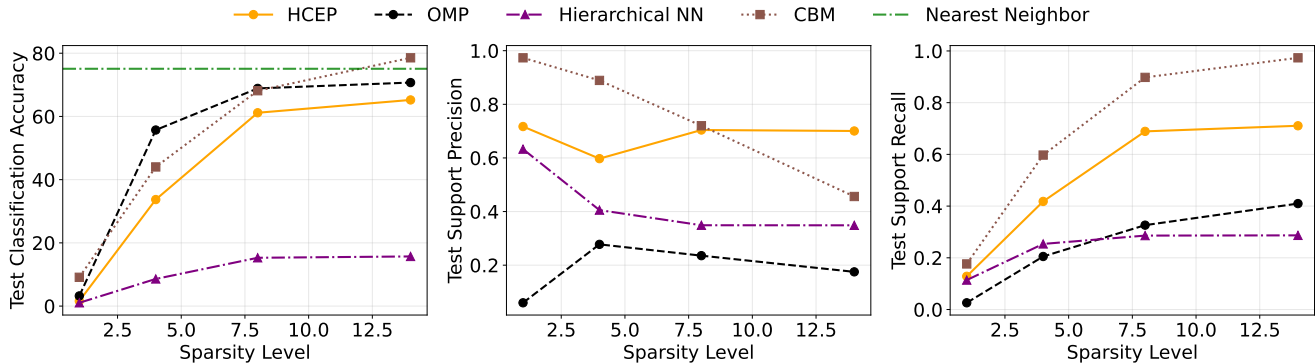


Figure 11. On ImageNet, HCEP achieves competitive accuracy while having higher concept precision/recall than sparse concept prediction baselines (OMP, Hierarchical NN).

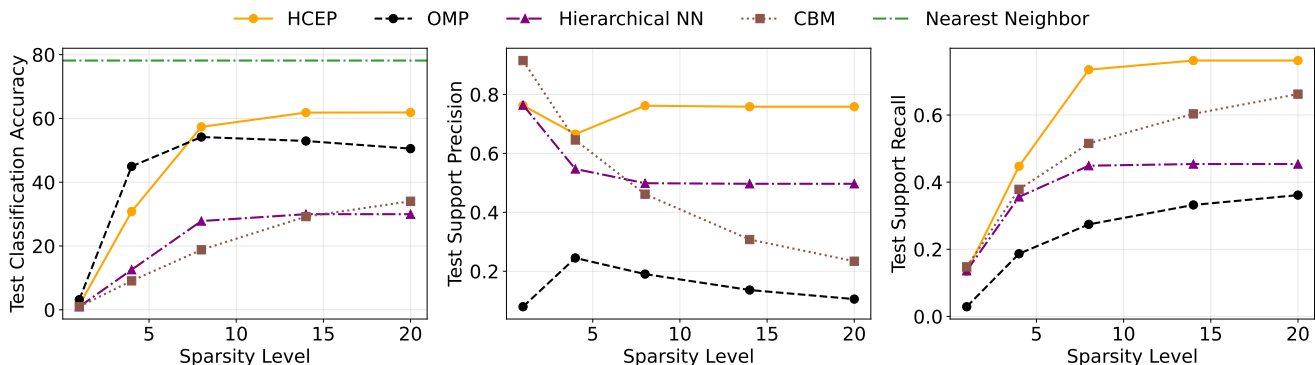


Figure 12. Interpretable image classification on ImageNet (12-shot) using SigLIP [72] embeddings. HCEP exhibits similar improvements in support precision and recall over baselines as with CLIP (cf. Fig. 11), demonstrating generality across vision-language models.

## F.2. Text Interpretation of Synset Differences

Optionally, to qualitatively evaluate alternative textual meanings of the synset differences, which form the atoms in our hierarchical dictionary, we use CLIP text embeddings and GPT-5 [51]. First, for each parent-child pair, we generate a text description of the difference between the parent and child synsets using GPT-5. Then, we pool the text embeddings of these descriptions to form a set of candidate concept embeddings. Next, for each synset difference vector, we find its top- $k$  neighbors among the candidate concept embeddings. Finally, we use GPT-5 to generate a summary of the top- $k$  neighbor descriptions, which helps interpret the synset difference. Table 1 shows example interpretations for several parent-child pairs in WordNet [47], the hierarchy underlying ImageNet.

## F.3. Additional Real-Data Experiment Details

**Model Architecture and Training Details.** We use CLIP-ViT-L/14 as the backbone. To train the linear classifier, we use the AdamW optimizer [41] with weight decay  $10^{-4}$  and learning rate  $10^{-1}$ . To train the CBM, we use the Adam op-

Table 1. Example text interpretations of child-parent synset differences produced via CLIP embeddings and GPT summarization.

Parent→Child Pair	Text Interpretation
bear → polar bear	thick matte white fur blending with snow.
container → basket	open-top woven or perforated sides with handles.
structure → lumbermill	vertical log-sawing machines and plank conveyors.
citrus → orange	round, bright orange, pebbled rind.

imizer with learning rate  $10^{-1}$  for 500 epochs. We provide the detailed hyperparameters in Table 2.

**Synset Difference Interpretations.** We use CLIP text embeddings [55] and GPT-5 [51] to generate textual interpretations of the synset differences. We provide the top-10 concepts for each parent-child pair in Table 3. We also include the GPT-5 prompt in Table 5.

**Ablation Study on Beam Size.** We perform an ablation

study on the beam size for Hierarchical OMP. We vary the beam size from 1 to 8 and evaluate concept recovery accuracy on ImageNette. We report the results in Fig. 14.

**Taxonomy Generation Prompt.** We use the taxonomy generation prompt from Zeng et al. [71] in Table 4 to generate the taxonomy for CIFAR-100.

Table 2. Key hyperparameters used in experiments for each dataset.

Hyperparameter	ImageNette	CIFAR-100	ImageNet
Batch size	4096	4096	16384
Classification training epochs	500	500	1000
HCEP beam size	8	16	32

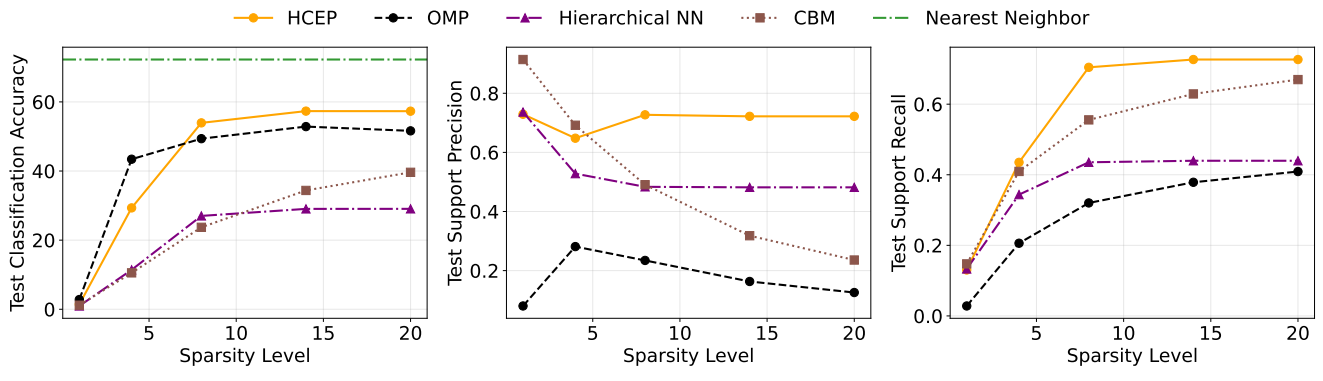


Figure 13. When we restrict the ImageNet training set to 25 images per class, HCEP outperforms all interpretable baselines.

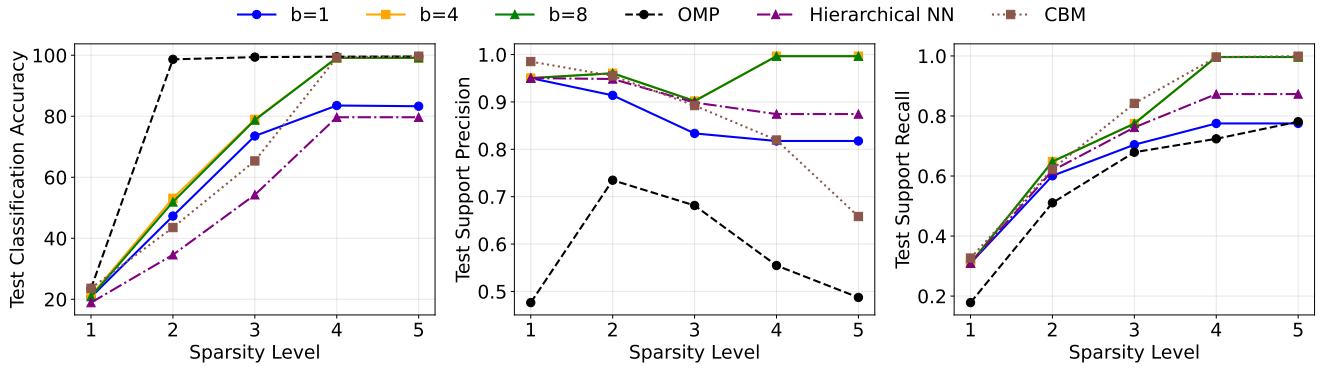


Figure 14. We vary the beam size over  $\{1, 4, 8\}$  and evaluate on ImageNette.

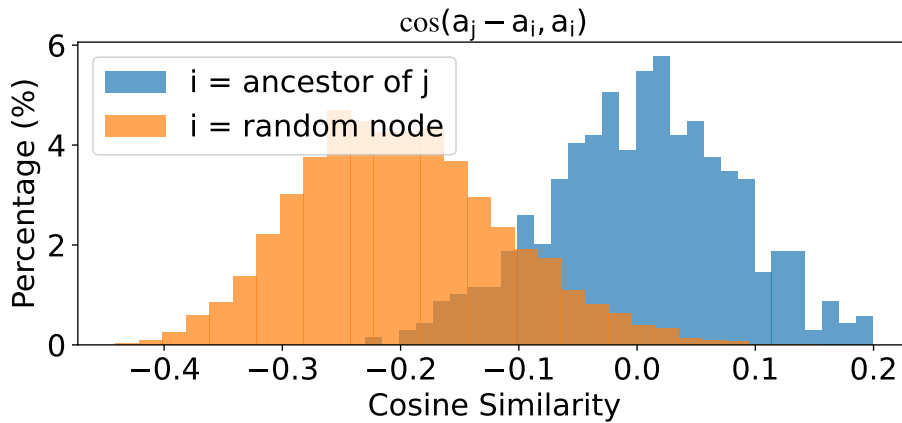


Figure 15. Observed hierarchical orthogonality on CIFAR-100. The cosine similarity between child-parent difference vectors and their parents is close to zero, while random non-parent pairs have significantly lower cosine similarity.

Table 3. Example text interpretations of child–parent synset differences produced via CLIP embeddings and GPT summarization, with the top-10 contributing concepts.

Parent→Child Pair	Top-10 Concepts	Text Interpretation
bear → polar bear	<ol style="list-style-type: none"> <li>1. Matte white texture</li> <li>2. Thick white winter fur coat</li> <li>3. White plumage in winter season</li> <li>4. White fur with cream patches</li> <li>5. Long, silky white coat</li> <li>6. White fur blending with surroundings</li> <li>7. Pure white fluffy coat</li> <li>8. White coat with lemon markings</li> <li>9. White cue ball</li> <li>10. Long, corded white coat</li> </ol>	thick matte white fur blending with snow.
container → basket	<ol style="list-style-type: none"> <li>1. Wicker baskets filled with baguettes</li> <li>2. Rectangular shopping baskets</li> <li>3. Stacked woven baskets</li> <li>4. Rear storage basket</li> <li>5. Rectangular open-top basket</li> <li>6. Woven rattan backrest</li> <li>7. Plastic shopping baskets</li> <li>8. Perforated cutlery basket in lower rack</li> <li>9. Woven rattan seating surfaces</li> <li>10. Rectangular basket frame</li> </ol>	open-top woven or perforated sides with handles.
structure → lumbermill	<ol style="list-style-type: none"> <li>1. Vertical log slicing machines</li> <li>2. Massive log cutting machines</li> <li>3. Metallic sawmill machinery</li> <li>4. Heavy-duty sawmill frames</li> <li>5. Conveyor belts with wood pieces</li> <li>6. Wooden board sorting systems</li> <li>7. Exposed wooden axles</li> <li>8. Wooden log feeding chutes</li> <li>9. Stacks of cut wooden planks</li> <li>10. Narrow wooden steering wheel</li> </ol>	vertical log-sawing machines and plank conveyors.

Table 4. LLM prompt used for taxonomy generation from root and leaf concepts.

### Taxonomy Generation Prompt

Given root concept `<root>` and leaf concepts `<leaves>`, generate a detailed hierarchical taxonomy that organizes these leaves under the root. Create multiple levels of intermediate category hierarchies to build a rich, fine-grained classification structure. Use as many hierarchical levels as needed to create meaningful semantic groupings and subgroupings.

The format is: 1. Parent Concept 1.1 Child Concept 1.1.1 Grandchild Concept.

**CRITICAL:** Every single leaf concept from the list must appear exactly as given in the taxonomy as the deepest level nodes. You may and should add multiple levels of intermediate concepts but do not add new leaf concepts. Before finishing, verify that each leaf concept from `<leaves>` appears in your taxonomy. Aim for depth and semantic richness in the hierarchy.

Table 5. LLM prompt used to generate child-vs-parent residual phrases.

### Residual Concept Generation Prompt

**TASK:** Generate a concise phrase (3-10 words) that describes what distinguishes "`<child>`" from its parent category "`<parent>`". This is for a hierarchical sparse model. A residual represents the visual difference between a child and parent category. Most correlated visual concepts (from CLIP embeddings): `<concepts_string>` **REQUIREMENTS:**

1. Generate ONE short phrase (3-10 words) that captures the key distinguishing features
2. Base your phrase on the correlated concepts provided above
3. Focus on the most salient visual features
4. Be specific and concrete, not vague or generic
5. Use natural language that a human would use to describe the difference
6. **IMPORTANT:** Do NOT use the synset names ("`<parent>`" or "`<child>`") in your phrase
7. **IMPORTANT:** Describe only the visual features, not the category name

**EXAMPLES OF GOOD PHRASES:** - "tawny coat with distinctive facial markings" (for a specific dog breed) - "long curved neck and pink coloration" (for flamingo vs bird) - "striped pattern and elongated body" (for a specific fish) - "metallic surface with cylindrical shape" (for a lighter)