NATIVE LOGICAL AND HIERARCHICAL REPRESENTATIONS WITH SUBSPACE EMBEDDINGS

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026027028

029

031

032

033

034

037

040

041

042

043

044

046 047

048

051

052

ABSTRACT

Traditional embeddings represent datapoints as vectors, which makes similarity easy to compute but limits how well they capture hierarchy, asymmetry and compositional reasoning. We propose a fundamentally different approach: representing concepts as learnable linear subspaces. By spanning multiple dimensions, subspaces can model broader concepts with higher-dimensional regions and nest more specific concepts within them. This geometry naturally captures generality through dimension, hierarchy through inclusion, and enables an emergent structure for logical composition, where conjunction, disjunction, and negation are mapped to linear operations. To make this paradigm trainable, we introduce a differentiable parameterization via soft projection matrices, allowing the effective dimension of each subspace to be learned end-to-end. We validate our approach on hierarchical and natural language inference benchmarks. Our method not only achieves state-of-the-art performance but also provides a more interpretable, geometrically-grounded model of entailment. Remarkably, the ability to perform logical composition with the learned concepts arises naturally from standard training objectives, without any direct supervision.

1 Introduction

Dense vector embeddings have become the bedrock of modern machine learning, underpinning systems from language models (LMs) (Devlin et al., 2019; Reimers & Gurevych, 2019) and vision-language models (VLMs) (Radford et al., 2021; Li et al., 2022), to retrieval augmented generation (RAG) systems (Lewis et al., 2020). By representing words, documents and images as points in high-dimensional space, these representations excel at capturing similarities in a scalable manner.

Despite their success, the efficacy of vector embeddings is limited by a geometric mismatch: the flat, symmetric structure of Euclidean space is ill-suited to the hierarchical and asymmetric nature of language and logic (Horn, 1972). Due to its symmetry, metrics like cosine similarity cannot capture directional relationships such as entailment or hyponymy; a high similarity between "dog" and "animal" fails to convey that one is a subtype of the other. Moreover, vector spaces lack native operators for logical conjunction and negation. This forces models to default to additive composition, effectively treating phrases as a bag-of-words. This explains why queries with negations often fail, with embeddings including the very concept meant for exclusion. Recent work confirms these flaws empirically, showing that even advanced models disregard logical connectives (Yuksekgonul et al., 2023; Moreira et al., 2025), requiring *ad-hoc* solutions (Weller et al., 2024; Gokhale et al., 2020; Zhang et al., 2025; Alhamoud et al., 2025). This inability to interpret nuanced instructions motivates our search for a framework that can natively represent these crucial relations.

We propose an alternative that extends Euclidean vector representations: instead of mapping a concept to a single vector, we embed it as a linear subspace of \mathbb{R}^d *i.e.*, the span of a set of basis vectors. This enables an interpretable geometric understanding of conceptual properties. First, generality and specificity are captured by the subspace dimension, with higher-dimensional subspaces denoting broader concepts *e.g.*, animal vs. dog. Secondly, hierarchy is naturally modeled by subspace inclusion, where a more specific concept's subspace is contained within a more general one. Finally, logical operations are directly mapped to linear-algebraic operations: conjunction as subspace intersection, disjunction as linear sum (span), and negation as the orthogonal complement (Fig. 1).

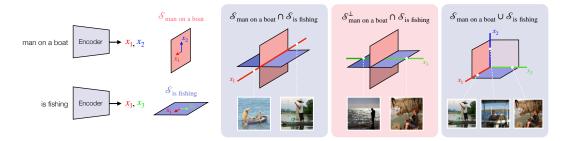


Figure 1: We embed concepts as linear subspaces of \mathbb{R}^d (left). These representations enable logical operations: subspace intersections e.g., "man on a boat" and "is fishing" (middle left); negation and composition e.g., orthogonal complement of "of man on a boat" and "is fishing" (middle right) and linear sums of subspaces, which yield a higher variance of instances (right).

A key challenge in learning subspaces is that dimensionality, or the number of basis vectors, is discrete and thus non-differentiable. Our technical contribution overcomes this by introducing a differentiable parameterization via soft projection matrices. Instead of selecting an integer dimension, we learn a set of vectors and modulate their individual importance, allowing each subspace to add or drop basis vectors as needed during training. Crucially, this approach remains grounded in Euclidean geometry, preserving full compatibility with standard training pipelines, Euclidean metrics and loss functions. This allows for seamless integration with highly efficient, dot-product-based search libraries Douze et al. (2025); Johnson et al. (2019), ensuring our method is scalable.

Beyond quantitative performance, our approach yields representations with emergent properties that are not explicitly optimized for. Remarkably, by training solely on entailment, our model learns embeddings that are inherently amenable to logical composition, supporting operations like conjunction, disjunction, and negation of queries. We also observe a strong correlation between the learned dimensionality of a subspace and the semantic generality of the concept it represents. This provides an interpretable measure of a concept's specificity that can be leveraged for compression.

We validate our framework across standard lexical and textual entailment benchmarks. Our method sets a new state of the art on WORDNET reconstruction, shows a stronger correlation with human judgments on HyperLex, and surpasses strong bi-encoder baselines on SNLI, demonstrating robust performance and generalization.

In summary, our key contributions are:

- A novel and differentiable method for learning subspace representations of language, featuring a data-dependent dimensionality that captures semantic specificity.
- An emergent structure for logical composition over natural language. We show that fundamental logical operations, such as conjunction, disjunction, and negation, arise naturally from standard training objectives, without any explicit logical supervision.
- A demonstration that these expressive representations remain tractable for large-scale retrieval by preserving compatibility with standard, highly optimized vector search pipelines.

2 Related Work

Most embedding methods, from Word2Vec (Mikolov et al., 2013) to multimodal models such as CLIP (Radford et al., 2021), rely on a simple idea: datapoints are represented as vectors in a high-dimensional metric space, where similarity is encoded by inner products or distances.

Limitations of Vector Representations. This prevalent vector-based view, while powerful for capturing co-occurrence patterns, exhibits limitations: the inner product cannot capture asymmetric relationships, such as entailment or hierarchies, without additional structural constraints or complex transformations. Recent empirical analyses have shed light on how language and vision-language encoder models represent hierarchies (Park et al., 2025; He et al., 2024) and logical constructs.

Remarkably, instead of capturing formal logical structure, vector embeddings behave akin to bagof-words representations (Yuksekgonul et al., 2023), failing to differentiate between positive and negated concepts (Gokhale et al., 2020; Singh et al., 2024; Moreira et al., 2025; Alhamoud et al., 2025). This limitation has motivated the creation of enhanced datasets and benchmarks with explicit negations (Quantmeyer et al., 2024; Weller et al., 2024; Zhang et al., 2025).

Hyperbolic Embeddings. Hyperbolic embeddings (Nickel & Kiela, 2018; 2017; Ganea et al., 2018a) exploit the exponential growth of hyperbolic space to model hierarchical structures compactly and encode transitive inclusion (Bai et al., 2021). Applications include hierarchical classification (Dhall et al., 2020), logical prediction (Xiong et al., 2022), and entailment reasoning (Poppi et al., 2025). However, they require complex Riemannian optimization, lack native logical reasoning, and struggle with non-hierarchical relations (Sala et al., 2018; Moreira et al., 2024).

Partial Order Embeddings. Partial order embeddings (Vendrov et al., 2016; Li et al., 2017) map entities into partially ordered spaces. Variants include positive operator embeddings (Lewis, 2019), probabilistic approaches such as Gaussians, mixtures, Beta distributions, box lattices (Vilnis & McCallum, 2015; Athiwaratkun & Wilson, 2018; Choudhary et al., 2021; Ren & Leskovec, 2020; Vilnis et al., 2018; Li et al., 2018; Ren et al., 2020), and entailment cones (Zhang et al., 2021; Pal et al., 2025; Ganea et al., 2018b; Yu et al., 2024). While effective for entailment, these methods typically lack principled logical operators, relying on heuristic approximations for disjunction and negation. Our subspace embeddings overcome this limitation. While subspace inclusion similarly models entailment, the key advantage is its algebraic closure: the intersection, sum, and orthogonal complement provide principled representations for conjunction, disjunction, and negation, respectively.

3 SUBSPACE REPRESENTATIONS

This paper presents a paradigm shift in embeddings: rather than representing a datapoint as a single vector $\boldsymbol{x} \in \mathbb{R}^d$, we represent it as a subspace $\mathcal{S} \subseteq \mathbb{R}^d$. To illustrate, consider Fig. 1. Instead of the traditional formulation, where the concept "man on a boat" is embedded as a single direction, we map it to \boldsymbol{x}_1 and \boldsymbol{x}_2 . Each vector encodes a variation of the underlying concept: \boldsymbol{x}_1 might represent a "man on a boat that is fishing" while \boldsymbol{x}_2 represents a "man on a boat that is not fishing". The concept "man on a boat" is then represented by the subspace $\mathcal{S}_{\text{man on a boat}} = \text{span}(\boldsymbol{x}_1, \boldsymbol{x}_2)$, encompassing all instances that align with either \boldsymbol{x}_1 , \boldsymbol{x}_2 , or any linear combination thereof, representing the space of all possible instances (Van Rijsbergen, 2004; Ganter & Wille, 2024).

Formally, we parameterize a subspace $\mathcal S$ as the span of $n \geq d$ learnable vectors $\boldsymbol X = [\boldsymbol x_1 \ \dots \ \boldsymbol x_n] \in \mathbb R^{d \times n}$. Let the thin singular value decomposition of $\boldsymbol X$ be $\boldsymbol U \boldsymbol \Sigma \boldsymbol V^\top$, with $\boldsymbol U \in \mathbb R^{d \times r}$ and $\boldsymbol U^\top \boldsymbol U = \boldsymbol I_r$. Then $\boldsymbol U$ is an orthonormal basis for the rank-r subspace $\mathcal S$. We can write an equivalent representation of $\mathcal S$ through its orthogonal projection operator,

$$P := X(X^{\top}X)^{\dagger}X^{\top} = UU^{\top} \in \mathbb{R}^{d \times d}, \tag{1}$$

where \dagger is the pseudoinverse. This projector is symmetric ($P^{\top} = P$), idempotent ($P^2 = P$), and its trace reveals the rank of S *i.e.*,

$$\operatorname{Tr}(\mathbf{P}) = \operatorname{Tr}(\mathbf{U}^{\top}\mathbf{U}) = \operatorname{Tr}(\mathbf{I}_r) = r.$$
 (2)

Subspace Similarity and Inclusion. Cosine similarity between vectors can be generalized to subspaces S_i , S_j , with orthonormal basis U_i , U_j , respectively, via their projection operators P_i and P_j ,

$$sim(\mathbf{P}_i, \mathbf{P}_j) := Tr(\mathbf{P}_i \mathbf{P}_j) = \|\mathbf{U}_i^{\mathsf{T}} \mathbf{U}_j\|_F^2 = \sum_{k=1}^m \cos^2(\theta_k),$$
(3)

where $\{\theta_k\}_{k=1}^m$ are the *principal angles* between \mathcal{S}_i and \mathcal{S}_j and $m = \min\{\operatorname{rank}(\mathcal{S}_i), \operatorname{rank}(\mathcal{S}_j)\}$. Each θ_k is the smallest possible angle between a unit vector in \mathcal{S}_i and a unit vector in \mathcal{S}_j , subject to orthogonality constraints on previously chosen directions. Thus, $\operatorname{sim}(P_i, P_j)$ measures the total squared alignment across the m most comparable directions of the two subspaces, or their degree of *overlap*. This recovers standard cosine similarity as a special case: if P_i and P_j are rank-one projectors onto unit vectors \mathbf{x}_i and \mathbf{x}_j , then $\operatorname{sim}(P_i, P_j) = (\mathbf{x}_i^{\top} \mathbf{x}_j)^2 = \cos^2(\angle(\mathbf{x}_i, \mathbf{x}_j))$.

An immediate consequence of Eqs. (2) and (3) is that we can quantify subspace inclusion via a normalized inclusion score (NIS) (Da Silva & Costeira, 2009):

$$NIS(\mathbf{P}_j \mid \mathbf{P}_i) := \frac{\sin(\mathbf{P}_i, \mathbf{P}_j)}{\operatorname{Tr}(\mathbf{P}_i)} \in [0, 1]. \tag{4}$$

This score attains 1 if and only if subspace i is contained within subspace j. This formulation allows for an intuitive interpretation as a Bayes-like conditional probability: the probability of an instance belonging to subspace j given it belongs to i.

3.1 ALGEBRAIC STRUCTURE OF SUBSPACES

The power of subspaces lies in their algebraic structure, which natively supports interpretable operations between concepts. Using projection operators lets us map logical relations such as conjunction (\land) , disjunction (\lor) and negation (\neg) into the subspace operations of intersection (\cap) , linear sum (+) and orthogonal complement (\bot) , respectively. These have tractable linear-algebraic representations, thus addressing the limitations of vector embeddings discussed in §1.

Conjunction $(i \wedge j)$. Corresponds to the intersection of subspaces $\mathcal{S}_{i \wedge j} = \mathcal{S}_i \cap \mathcal{S}_j$. Any vector in $\mathcal{S}_{i \wedge j}$ is an element of \mathcal{S}_i and \mathcal{S}_j . The product $P_i P_j$ is an orthogonal projection onto $\mathcal{S}_i \cap \mathcal{S}_j$ if and only if P_i and P_j commute. In the general case, $P_{i \wedge j} = \lim_{n \to \infty} (P_i P_j)^n$. In Fig. 1, the intersection $\mathcal{S}_{\text{man on a boat}} = \operatorname{span}(\boldsymbol{x}_1, \boldsymbol{x}_2)$ and $\mathcal{S}_{\text{is fishing}} = \operatorname{span}(\boldsymbol{x}_1, \boldsymbol{x}_3)$ yields $\mathcal{S}_{\text{man fishing on a boat}} = \operatorname{span}(\boldsymbol{x}_1)$.

Disjunction $(i \lor j)$. Corresponds to the span (linear sum) of subspaces: $S_{i \lor j} = S_i + S_j$. Any vector in $S_{i \lor j}$ is a linear combination of elements in S_i or in S_j . For commuting subspaces, the projection onto $S_{i \lor j}$ satisfies $P_{i \lor j} = P_i + P_j - P_{i \land j}$. In Fig. 1, the linear sum $S_{\text{man fishing on a boat}} = \text{span}(\boldsymbol{x}_1)$ and $S_{\text{man fishing not on a boat}} = \text{span}(\boldsymbol{x}_3)$ yields $S_{\text{man fishing}} = \text{span}(\boldsymbol{x}_1, \boldsymbol{x}_3)$.

Complement $(\neg i)$. Corresponds to the subspace of all vectors orthogonal to the subspace: $\mathcal{S}_{\neg i} = \mathcal{S}_i^{\perp}$. The projection operator onto \mathcal{S}_i^{\perp} is given by $\mathbf{P}_{\neg i} = \mathbf{I} - \mathbf{P}_i$. In Fig. 1, the complement of $\mathcal{S}_{\text{man on a boat}} = \text{span}(\mathbf{x}_1, \mathbf{x}_2)$ is given by $\mathcal{S}_{\text{man on a boat}}^{\perp} = \text{span}(\mathbf{x}_3)$.

3.2 Representing Subspaces as Soft Projection Operators

While the orthogonal projector from Eq. (1) offers a rich and interpretable parameterization of a subspace, its optimization poses a challenge for gradient-based methods. Since the rank of a subspace is integer-valued, the space of all subspaces (a union of Grassmannian manifolds) is stratified and non-differentiable across rank changes. This makes it hard to simultaneously learn orientation and dimensionality via gradient-based methods.

Soft Projection Operators. To overcome the challenges associated with learning adaptive-rank subspaces we introduce a relaxation of the projection operator in Eq. (1). For a rank-r subspace \mathcal{S} spanned by the columns of $\mathbf{X} \stackrel{\text{SVD}}{=} \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top} \in \mathbb{R}^{d \times n}$, where $\mathbf{\Sigma} = \text{diag}(\{\sigma_i\}_{i=1}^r)$ and $\mathbf{U} \in \mathbb{R}^{d \times r}$ is the orthonormal basis of \mathcal{S} , we define a soft projector via Tikhonov regularization

$$\tilde{\boldsymbol{P}} := \boldsymbol{X} (\boldsymbol{X}^{\top} \boldsymbol{X} + \lambda \boldsymbol{I})^{-1} \boldsymbol{X}^{\top} = \boldsymbol{U} \operatorname{diag} \left(\left\{ \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \right\}_{i=1}^r \right) \boldsymbol{U}^{\top}, \quad \lambda > 0.$$
 (5)

Unlike a true projector ($P^2 = P$), \tilde{P} is a *soft projector*: its eigenvalues vary smoothly in [0,1) rather than being binary. This makes the operator differentiable with respect to both orientation and rank, avoiding hard rank jumps and enabling gradual changes in dimensionality. Geometrically, this relaxation replaces the stratified manifold of projectors with a smooth manifold of PSD operators. From a Bayesian point of view, it corresponds to a Gaussian prior with precision λI .

For small values of λ , the soft projectors in Eq. (5) provide accurate surrogates for the algebraic operations and metrics introduced in §3.1. The approximation error depends primarily on the weakest nonzero singular value σ_r of X, being upper bounded by (see Appendix A)

$$\epsilon(\sigma_r, \lambda) = \lambda/(\sigma_r^2 + \lambda).$$
 (6)

Table 1: Soft approximations of projection operations derived from X_i and X_j . Errors are in operator norm, except rank (relative absolute error). σ_r , η_r are the weakest non-null singular values of X_i and X_j . $\epsilon(\sigma_r, \lambda) = \lambda/(\sigma_r^2 + \lambda)$.

	Projector	Negation	Intersection	Linear sum	Rank
Exact	$oldsymbol{X}(oldsymbol{X}oldsymbol{X}^ op)^\daggeroldsymbol{X}^ op$	I - P	$oldsymbol{P}_i oldsymbol{P}_j$	$oldsymbol{P}_i + oldsymbol{P}_j - oldsymbol{P}_i oldsymbol{P}_j$	$\operatorname{Tr}(oldsymbol{P})$
Soft	$\boldsymbol{X}(\boldsymbol{X}\boldsymbol{X}^{\top} + \lambda \boldsymbol{I})^{-1}\boldsymbol{X}^{\top}$	I - P	$m{P}_im{P}_j$	$ ilde{m{P}}_i + ilde{m{P}}_j - ilde{m{P}}_i ilde{m{P}}_j$	$\mathrm{Tr}(oldsymbol{P})$
Error	$\epsilon(\sigma_r,\lambda)$	$\epsilon(\sigma_r,\lambda)$	$\epsilon(\sigma_r,\lambda) + \epsilon(\eta_r,\lambda)$	$2(\epsilon(\sigma_r,\lambda)+\epsilon(\eta_r,\lambda))$	$\epsilon(\sigma_r,\lambda)$

As $\lambda \to 0$, $\epsilon(\sigma_r, \lambda) \to 0$ and we recover the orthogonal projection operator $\tilde{P} \to P$, while larger λ enforces smoother, more regularized projectors. Table 1 summarizes how each operation is approximated using \tilde{P} and the resulting deviation from the orthogonal projector ($\lambda = 0$).

Subspace Projection Head (SPH). To bridge our subspace representations with transformer models, we introduce the *Subspace Projection Head (SPH)*. A transformer first encodes text inputs into a contextualized hidden state $\boldsymbol{H} \in \mathbb{R}^{h \times m}$ (where m is sequence length, h is hidden dimension). The SPH transforms this hidden state \boldsymbol{H} into a fixed-size set of n vectors $\boldsymbol{X} \in \mathbb{R}^{d \times n}$ that span the subspace \mathcal{S} , and computes the corresponding soft projector $\tilde{\boldsymbol{P}}$.

We map the hidden state H into a sequence-length-invariant subspace in three stages. First, we augment the transformer with a set of n learnable query vectors $Q \in \mathbb{R}^{h \times n}$. These queries, attend to H (acting as keys and values) via Multi-Head Attention (MHA), pooling n embeddings X',

$$X' = MHA(query = Q, key = H, value = H) \in \mathbb{R}^{h \times n}.$$
 (7)

This ensures invariance to sequence length m. However, the rank of X' is still limited: since each head outputs a linear combination of the columns of H, then $\mathrm{rank}(X') \leq n_{\mathrm{heads}} \cdot \mathrm{rank}(H) \leq m \cdot n_{\mathrm{heads}}$. We address this via a Multi-Layer Perceptron (MLP) which maps the n h-dimensional vectors from the MHA output to \mathbb{R}^d as $X = \mathrm{MLP}(X')$. This yields the subspace matrix $X \in \mathbb{R}^{d \times n}$, which spans the subspace. Finally, the soft projector \tilde{P} is computed from X using Eq. (5).

3.3 Training Methodology

We learn subspaces end-to-end via gradient descent, requiring no special pretraining, or training constraints. Depending on the downstream task, we employ one of the following loss functions.

Reconstruction. For similarity-based tasks, we use an InfoNCE loss (van den Oord et al., 2019) with the subspace similarity computed via $sim(\tilde{P}_i, \tilde{P}_j)$, from Eq. (3).

Link Prediction. In link prediction tasks, we optimize the normalized inclusion score NIS($\vec{P}_i \mid \tilde{P}_j$) from Eq. (4) directly and consider the margin loss (Vendrov et al., 2016)

$$L = \sum_{i,j \in \mathcal{P}} [\gamma_{+} - \text{NIS}(\tilde{\mathbf{P}}_{i} \mid \tilde{\mathbf{P}}_{j})]_{+} + \sum_{i,j \in \mathcal{N}} [\text{NIS}(\tilde{\mathbf{P}}_{i} \mid \tilde{\mathbf{P}}_{j}) - \gamma_{-}]_{+}, \tag{8}$$

where $[\cdot]_+$ denotes the ReLU function. Here, $\gamma_+, \gamma_- \in (0,1)$ are the positive and negative margins and \mathcal{P} and \mathcal{N} the set of positives and negatives, respectively.

NLI Classification. Textual Entailment presents a unique challenge, requiring not just a measure of inclusion but also an explicit model of neutrality. For a premise p and hypothesis h, we model the relation $Y \in \{E, N, C\}$ (entailment, neutral, contradiction) as a discrete latent variable. For $Y \in \{E, C\}$, we assume the generative process for $S = \text{NIS}(\tilde{P}_h \mid \tilde{P}_p)$

$$S \mid (Y = y) \sim \text{Beta}(\alpha_y, \beta_y), \quad y \in \{E, C\},$$
 (9)

with $\alpha_y \leq \beta_y$ if y = C and $\beta_y \leq \alpha_y$ if y = E. For neutrals, subspace inclusion does not provide a reliable signal. Instead, we model neutrality independently by an MLP as

$$P(Y = y \mid \tilde{\mathbf{P}}_{p}, \tilde{\mathbf{P}}_{h}) := \sigma\left(\text{MLP}\left(\tilde{\mathbf{P}}_{p}, \tilde{\mathbf{P}}_{h}, \tilde{\mathbf{P}}_{p}\tilde{\mathbf{P}}_{h}, \tilde{\mathbf{P}}_{h}\tilde{\mathbf{P}}_{p}\right)\right), \quad y = N$$
(10)

Table 2: WORDNET **reconstruction**. mAP = Mean Average Precision, MR = Mean Rank, ρ = Spearman correlation between taxonomy rank and subspace dimension or norm (for \mathcal{P}^{10} , \mathcal{H}^{10}).

Method		Nouns		Verbs			
17100110u	mAP (†)	MR (↓)	$\rho \left(\uparrow \right)$	mAP (†)	MR (↓)	$\rho \left(\uparrow \right)$	
Euclidean (\mathbb{R}^{128})	95.1	1.31	_	98.6	1.04	_	
Poincaré (\mathcal{P}^{10})	86.5	4.02	58.5	91.2	1.35	55.1	
Lorentz (\mathcal{H}^{10})	92.8	2.95	59.5	93.3	1.23	56.6	
Subspaces (SE ¹²⁸)	98.6	1.04	68.0	99.9	1.00	67.0	

where $\sigma(\cdot)$ denotes the sigmoid function. Assuming uniform priors for entailment and contradiction classes, conditional on non-neutrality, we compute posterior probabilities for y=E and y=C, denoted $P(Y=y\mid S=s,Y\in\{E,C\})$. The final posterior probabilities for $y\in\{E,C\}$ are then derived by combining the MLP output for neutrality with the Beta posteriors for non-neutrality:

$$P(Y = y \mid \tilde{P}_p, \tilde{P}_h, S = s) = (1 - P(Y = N \mid \tilde{P}_p, \tilde{P}_h))P(Y = y \mid S = s, Y \neq N),$$
 (11)

for $y \in \{E, C\}$. The posteriors in Eqs. (10) and (11) are optimized via a cross-entropy loss.

A key insight into how these losses shape the subspaces is revealed by the gradient dynamics. As derived in Appendix B, $\nabla_{X_i} \text{sim}(\tilde{P}_i, \tilde{P}_j)$ encourages subspace i to expand along the principal directions of subspace j that it currently lacks. This update naturally promotes subspace inclusion, and the gradient vanishes once one subspace is contained within the other, leading to stable convergence.

Efficiency Considerations. While computing \tilde{P} from $X \in \mathbb{R}^{d \times n}$ is $\mathcal{O}(n^3)$ in the number of vectors n, and has a memory footprint that scales with d^2 , where d is the ambient dimension, our approach is practical for two key reasons. First, the model learns a data-dependent rank for each subspace. As our experiments demonstrate, this allows for considerable compression of \tilde{P} via low-rank approximations. Second, the subspace similarity (3) and NIS (4) are equivalent to dot products between the vectorized matrices: $\mathrm{Tr}(\tilde{P}_i\tilde{P}_j) = \mathrm{vec}(\tilde{P}_i)^{\top}\mathrm{vec}(\tilde{P}_j)$. This allows our subspaces to be indexed by highly optimized vector search libraries, making large-scale retrieval feasible.

4 EXPERIMENTS

We empirically validate our embeddings' ability to model large-scale hierarchies and textual entailment on a suite of benchmarks including WORDNET (Miller, 1995) reconstruction in §4.1 and link prediction in §4.2, HyperLex (Vulić et al., 2017), and SNLI (Bowman et al., 2015) in §4.3.

4.1 WORDNET RECONSTRUCTION

In WORDNET's reconstruction task, all edges from the full transitive closure of the noun and verb hypernymy hierarchies are used for training and testing. The goal is to assess the capacity of the representations to capture known hierarchical relations by providing only pairwise relations.

Experimental Details. Each node in the graph is represented by a soft projection matrix (5), with $\lambda=0.2$, parameterized by a matrix $\boldsymbol{X}_i\in\mathbb{R}^{128\times128}$. For each training edge (u,v), we sample 19 nodes $v'\neq u$ such that neither (u,v') nor (v',u) are in the train split and optimize InfoNCE using Adam (Kingma & Ba, 2017). During evaluation, we first compute the subspace similarity $\mathrm{Tr}(\tilde{\boldsymbol{P}}_u\tilde{\boldsymbol{P}}_v)$ of every edge (u,v) in the transitive closure. We then rank each of these scores among those of all node pairs that are not connected in the transitive closure. Based on these rankings, we report the mean rank (MR) and the mean average precision (mAP). Additional details in Appendix C.1.

Reconstruction Results. Our method achieves state-of-the-art performance on the WordNet reconstruction. As shown in Table 2, our subspace representations (SE¹²⁸) significantly outperform both Hyperbolic (Poincaré \mathcal{P}^{10} and Lorentz \mathcal{H}^{10} models), and Euclidean embeddings (\mathbb{R}^{128}) baselines, with a near-perfect reconstruction on the shallower verb hierarchy.

Table 3: HYPERLEX lexical entailment Spearman's rank correlation (WORDNET embeddings).

	\mathbb{R}^5	\mathcal{P}^5	DOE-A ⁵⁰	$SE^{128} (\lambda = 0.2)$	$SE^{128} (\lambda = 0.6)$
$\rho (\uparrow)$	0.389	0.512	0.590	0.683	0.734

Table 4: WORDNET noun **link prediction** F1-Scores (†). Superscript denotes dimension.

% Non-Basic Edges	\mathbb{R}^{10}	OE^{10}	\mathcal{P}^{10}	Cones ¹⁰	Disk ¹⁰	UHS ¹⁰	SE^{64}	SE^{128}
0%	29.4	43.0	29.0	32.4	36.5	52.2	48.9	53.4
10%	75.4	69.7	71.5	84.9	79.5	89.4	93.7	94.2
25%	78.4	79.4	82.1	90.8	90.5	95.7	95.7	96.0
50%	78.1	84.1	85.4	93.8	94.2	97.0	95.9	95.4

To assess how these representations generalize to graded lexical entailment, we evaluated them on the HYPERLEX noun subset without fine-tuning (see Appendix C.3). We quantify entailment using the NIS from Eq. (4), selecting the synset pair with maximal similarity for disambiguation (Athiwaratkun & Wilson, 2018). As reported in Table 3, our embeddings demonstrate a significantly stronger correlation with human judgments than prior work. Our approach achieves a Spearman's ρ of 0.73 ($\lambda=0.6$), substantially outperforming Poincaré and Gaussian embedding baselines.

4.2 WORDNET LINK PREDICTION

In the link prediction task, we evaluate generalization from sparse supervision. We split the set of edges from the transitive closure that are not part the original graph (non-basic edges) into train (90%), validation (5%) and test (5%) using the data split from Suzuki et al. (2019).

Experimental Details. To assess how the percentage of the transitive closure seen during training impacts performance, we created partial training edge coverages by randomly sampling 0%, 10%, 25% or 50% of non-basic edges, to which we append all the basic edges. We considered two ambient space dimensions d=64 and d=128, setting the number of vectors as n=d in each case. Training was performed by optimizing the margin loss defined in Eq. (8). During evaluation, for each positive test edge, we consider 10 negative test edges: half with a corrupted head, and half with a corrupted tail. We classify edges by thresholding the NIS from Eq. (4) and report the classification F1-Score.

Link Prediction Results. Link prediction results are shown in Table 4. We compare against Euclidean embeddings (\mathbb{R}^{10}), Order Embeddings ($\mathrm{OE^{10}}$), Poincaré (\mathcal{P}^{10}) Nickel & Kiela (2017), Hyperbolic Entailment Cones (Cones¹⁰) (Ganea et al., 2018b), Hyperbolic Disk Embeddings (Disk¹⁰) (Suzuki et al., 2019) and the Umbral Half-Space embeddings (UHS¹⁰) (Yu et al., 2024). Subspace embeddings SE⁶⁴ and SE¹²⁸ outperform the baselines across most supervision levels. SE¹²⁸, in particular, offers a considerable improvement when training with sparser supervision. This underscores the ability of subspace representations to infer hierarchical relations even from weak supervision.

4.3 SNLI

We conducted experiments on NLI using the SNLI dataset. SNLI comprises 550,152 training, and 10,000 validation/test premise (p) - hypothesis (h) pairs, each annotated with one of three labels: entailment, neutral, or contradiction. We consider two regimes: 3-way, and 2-way classification (entailment vs non-entailment). For a fair comparison, we benchmarked bi-encoder baselines, using the all-miniLM-L6-v2 and mpnet-base-v2 models with a shallow MLP classifier. We considered two variants: MLP(p,h), using concatenated premise p and hypothesis h embeddings, and MLP(p,h,p-h). In our models, we map p and h to soft projectors via our SPH module ($\lambda=0.05$). All models were trained with a cross-entropy loss. Additional details are provided in Appendix D.

Results. As shown in Table 5, our approach consistently outperforms bi-encoder baselines. For reference, we also include two GRU-based hierarchical approaches: Order Embeddings and Hyperbolic Neural Networks (HNN) Ganea et al. (2018a), which do not model neutrality. Crucially, in the 2-way setting, our method, which relies solely on subspace inclusion, consistently outperforms the MLP baselines, with a more interpretable mechanism.

Table 5: SNLI test accuracy: 2-way (entailment vs non-entailment) and 3-way (+Neutral).

Method	2-way	3-way				
Order Embeddings (GRU)	88.60	_				
HNN (GRU)	81.19	_				
all-miniLM-L6-v2 (22.7m p	arams)					
MLP(p, h)	90.48	83.63				
MLP(p, h, p - h)	91.06	84.74				
SPH (SE^{64})	91.02	84.62				
SPH (SE ¹²⁸)	91.12	85.24				
mpnet-base-v2 (109m params)						
MLP(p, h)	90.86	83.77				
MLP(p, h, p - h)	91.63	86.14				
SPH (SE^{64})	92.27	85.66				
SPH (SE ¹²⁸)	92.21	86.50				

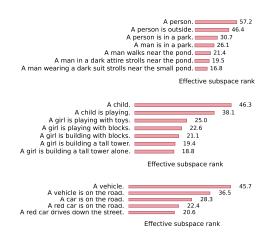


Figure 2: Example effective subspace ranks.

5 QUALITATIVE ANALYSIS

A key finding of our work is that our framework learns an interpretable geometry that maps the hierarchical structure of language onto the representations. We confirm this empirically on SNLI, using our SE¹²⁸ embeddings. In Fig. 3, we plot the histogram of the NIS (4) for premise-hypothesis pairs encoded with out mpnet-base-v2 (SE¹²⁸) subspace model. We observe that, for entailment this metric is concentrated towards 1, for contradictions it skews towards 0, and for neutrals it is centered around 0.5. This confirms that the NIS reflects the underlying entailment structure via subspace inclusion: each premise subspace is contained within the hypotheses subspaces it entails.

Rank as a Measure of Generality. A direct consequence of this is that a subspace's effective rank, as measured by $\mathrm{Tr}(\tilde{P})$, becomes an emergent measure of semantic generality. For a specific concept to be nested within many broader ones, it must occupy a lower-dimensional subspace. This property is confirmed quantitatively by the high Spearman correlation (ρ) between WORDNET nouns' true hierarchical positions (distance from root) and their learned effective rank in Table 2. We provide additional visual confirmation of this principle. Fig. 4 shows how the effective rank of WORDNET's nouns grows with the number of

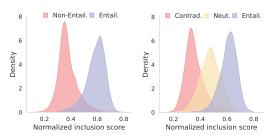


Figure 3: Histogram of the NIS for SNLI's test premise-hypothesis pairs encoded with our SE¹²⁸ model. Left: Two-way; Right: Three-way.

their descendants. The annotated chain from the specific *homo sapiens* to the root noun *entity* clearly illustrates this monotonic increase. Fig. 2 shows the same phenomenon for three entailment sequences. We observe again that, the *effective rank* of each sentence increases as we go from a specific description to general one e.g., "A car is on the road." \rightarrow "A vehicle."

Dimensionality Reduction. This learned structure, where the rank encodes specificity, makes our embeddings inherently compressible. Since each subspace dynamically allocates the dimensions needed to represent each concept, we can perform post-training compression via truncated SVD, with minimal performance loss. To assess this capability, we approximated WORDNET and SNLI embeddings \tilde{P}_i by retaining singular values greater than a threshold $\tau \in [0,1]$ and plot the reconstruction mAP, in the case of WordNet, or the two-way accuracy, for SNLI, as well as the average subspace rank, as a function τ . As shown in Fig. 5, the learned subspaces exhibit rapid spectral decay in both experiments, allowing for compression of up to $4\times$ with negligible impact on task performance. This paves the way for a new class of embeddings where representational complexity is not fixed, but a learned, data-driven property.

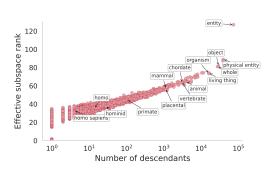


Figure 4: *Effective rank* $\operatorname{Tr}(\tilde{\boldsymbol{P}})$ vs number of descendants of WORDNET nouns.

Emergent Logical Reasoning. A key advantage of subspace embeddings is their emergent logical compositionality, which arises directly from the geometry of the embeddings, without explicit training signals. Fig. 6 provides an example illustrating this inherent compositionality, for conjunctions $\hat{P}_i\hat{P}_j$ and negations $I - \hat{P}$, in a retrieval setting. For a query formed by a logical combination of concept subspaces, we retrieve images from Flickr30k (Young et al., 2014) whose caption subspaces have the largest $NIS(\tilde{P}_{query} \mid \tilde{P}_{caption})$. Each caption subspace is computed with our mpnet-base-v2 + SPH (SE¹²⁸) model, fine-tuned on SNLI. The results demonstrate that subspaces enable compositional retrieval, allowing for the search of novel concepts or query editing through geometric operations. Additional examples in Appendix E.

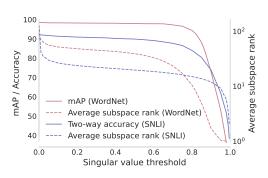


Figure 5: Accuracy, mAP and average rank as a function of the singular value threshold.



Figure 6: Flickr30k retrieval from composition of natural language queries.

6 EFFICIENCY ANALYSIS

By vectorizing the similarity (3) and the NIS (4), we can make our embeddings compatible with fast search libraries like FAISS. This contrasts with non-Euclidean embeddings requiring brute-force search. We benchmarked retrieval latency on CPU over the 155,070 Flickr30k captions (batch-size 128). The results in Table 6 show that SE¹²⁸ is nearly $8 \times$ **faster** than a 10D Poincaré (\mathcal{P}^{10}) baseline. The encoding overhead introduced by the SPH is also minimal, averaging at an additional 0.12ms/query on GPU (Appendix F).

Table 6: Search efficiency.

Latency (ms/query)			
$\overline{\mathcal{P}^{10}}$	3.64 ± 0.13		
SE^{128}	0.47 ± 0.02		

7 CONCLUSION

This paper introduced *subspace embeddings*, a novel paradigm that addresses the limitations of vector representations in capturing logical structure and asymmetric, or hierarchical, relations. By representing concepts as subspaces, our framework naturally encodes generality through dimensionality and hierarchy through inclusion. Our evaluation across hierarchical and entailment tasks reveals the power of this inductive bias: it not only achieves state-of-the-art results but also gives rise to an emergent structure for logical composition without explicit supervision. The linearity of the core operations and metrics ensures compatibility with efficient vector search pipelines. Overall, our results establish subspace embeddings as a bridge between representation learning and logical reasoning, opening avenues for new representations that exploit the structural nature of data.

REFERENCES

- Kumail Alhamoud, Shaden Alshammari, Yonglong Tian, Guohao Li, Philip HS Torr, Yoon Kim, and Marzyeh Ghassemi. Vision-language models do not understand negation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 29612–29622, 2025.
- Ben Athiwaratkun and Andrew Gordon Wilson. Hierarchical density order embeddings. In 6th International Conference on Learning Representations, ICLR 2018, 2018.
- Yushi Bai, Zhitao Ying, Hongyu Ren, and Jure Leskovec. Modeling heterogeneous hierarchies with relation-specific hyperbolic cones. volume 34, pp. 12316–12327, 2021.
 - Samuel Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 632–642, 2015.
 - Nurendra Choudhary, Nikhil Rao, Sumeet Katariya, Karthik Subbian, and Chandan Reddy. Probabilistic entity representation model for reasoning over knowledge graphs. volume 34, pp. 23440–23451, 2021.
 - Nuno Pinho Da Silva and Joao Paulo Costeira. The normalized subspace inclusion: Robust clustering of motion subspaces. In *2009 IEEE 12th International Conference on Computer Vision*, pp. 1444–1450. IEEE, 2009.
 - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.
 - Ankit Dhall, Anastasia Makarova, Octavian Ganea, Dario Pavllo, Michael Greeff, and Andreas Krause. Hierarchical image classification using entailment cone embeddings. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 836–837, 2020.
 - Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library, 2025. URL https://arxiv.org/abs/2401.08281.
 - Octavian Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic neural networks. volume 31, 2018a.
 - Octavian Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic entailment cones for learning hierarchical embeddings. In *International conference on machine learning*, pp. 1646–1655. PMLR, 2018b.
- Bernhard Ganter and Rudolf Wille. Formal concept analysis: mathematical foundations. Springer Nature, 2024.
- Tejas Gokhale, Pratyay Banerjee, Chitta Baral, and Yezhou Yang. Vqa-lol: Visual question answering under the lens of logic. In *European conference on computer vision*, pp. 379–396. Springer, 2020.
- Yuan He, Moy Yuan, Jiaoyan Chen, and Ian Horrocks. Language models as hierarchy encoders. volume 37, pp. 14690–14711, 2024.
- Laurence Robert Horn. On the semantic properties of logical operators in English. University of California, Los Angeles, 1972.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.
 - Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL https://arxiv.org/abs/1412.6980.

- Martha Lewis. Compositional hyponymy with positive operators. In *Proceedings of the Inter*national Conference on Recent Advances in Natural Language Processing (RANLP 2019), pp. 638–647, 2019.
 - Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33: 9459–9474, 2020.
 - Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.
 - Xiang Li, Luke Vilnis, and Andrew McCallum. Improved representation learning for predicting commonsense ontologies, 2017. URL https://arxiv.org/abs/1708.00549.
 - Xiang Li, Luke Vilnis, Dongxu Zhang, Michael Boratko, and Andrew McCallum. Smoothing the geometry of probabilistic box embeddings. In *International Conference on Learning Representations*, 2018.
 - Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. volume 26, 2013.
 - George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11): 39–41, 1995.
 - Gabriel Moreira, Manuel Marques, João Paulo Costeira, and Alexander Hauptmann. Hyperbolic vs euclidean embeddings in few-shot learning: Two sides of the same coin. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2082–2090, 2024.
 - Gabriel Moreira, Alexander Hauptmann, Manuel Marques, and João Paulo Costeira. Learning visual-semantic subspace representations. 2025.
 - Maximillian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. volume 30, 2017.
 - Maximillian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz model of hyperbolic geometry. In *International conference on machine learning*, pp. 3779–3788. PMLR, 2018.
 - Avik Pal, Max van Spengler, Guido Maria D'Amely di Melendugno, Alessandro Flaborea, Fabio Galasso, and Pascal Mettes. Compositional entailment learning for hyperbolic vision-language models. In *The Thirteenth International Conference on Learning Representations*, 2025.
 - Kiho Park, Yo Joong Choe, Yibo Jiang, and Victor Veitch. The geometry of categorical and hierarchical concepts in large language models, 2025. URL https://arxiv.org/abs/2406.01506.
 - Tobia Poppi, Tejaswi Kasarla, Pascal Mettes, Lorenzo Baraldi, and Rita Cucchiara. Hyperbolic safety-aware vision-language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 4222–4232, 2025.
 - Vincent Quantmeyer, Pablo Mosteiro, and Albert Gatt. How and where does clip process negation? In *Proceedings of the 3rd Workshop on Advances in Language and Vision Research (ALVR)*, pp. 59–72, 2024.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bertnetworks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3982–3992, 2019.

- H Ren, W Hu, and J Leskovec. Query2box: Reasoning over knowledge graphs in vector space using box embeddings. In *International Conference on Learning Representations (ICLR)*, 2020.
 - Hongyu Ren and Jure Leskovec. Beta embeddings for multi-hop logical reasoning in knowledge graphs. volume 33, pp. 19716–19726, 2020.
 - Frederic Sala, Chris De Sa, Albert Gu, and Christopher Ré. Representation tradeoffs for hyperbolic embeddings. In *International conference on machine learning*, pp. 4460–4469. PMLR, 2018.
 - Jaisidh Singh, Ishaan Shrivastava, Mayank Vatsa, Richa Singh, and Aparna Bharati. Learn "no" to say "yes" better: Improving vision-language models via negations, 2024. URL https://arxiv.org/abs/2403.20312.
 - Ryota Suzuki, Ryusuke Takahama, and Shun Onoda. Hyperbolic disk embeddings for directed acyclic graphs. In *International Conference on Machine Learning*, pp. 6066–6075. PMLR, 2019.
 - Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2019. URL https://arxiv.org/abs/1807.03748.
 - Cornelis Joost Van Rijsbergen. *The geometry of information retrieval*. Cambridge University Press, 2004.
 - Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. Order-embeddings of images and language, 2016. URL https://arxiv.org/abs/1511.06361.
 - Luke Vilnis and Andrew McCallum. Word representations via gaussian embedding, 2015. URL https://arxiv.org/abs/1412.6623.
 - Luke Vilnis, Xiang Li, Shikhar Murty, and Andrew Mccallum. Probabilistic embedding of knowledge graphs with box lattice measures. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 263–272, 2018.
 - Ivan Vulić, Daniela Gerz, Douwe Kiela, Felix Hill, and Anna Korhonen. HyperLex: A large-scale evaluation of graded lexical entailment. *Computational Linguistics*, 43(4):781–835, December 2017. doi: 10.1162/COLLa_00301. URL https://aclanthology.org/J17-4004/.
 - Orion Weller, Dawn Lawrie, and Benjamin Van Durme. Nevir: Negation in neural information retrieval. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2274–2287, 2024.
 - Bo Xiong, Michael Cochez, Mojtaba Nayyeri, and Steffen Staab. Hyperbolic embedding inference for structured multi-label prediction. volume 35, pp. 33016–33028, 2022.
 - Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014.
 - Tao Yu, Toni JB Liu, Albert Tseng, and Christopher De Sa. Shadow cones: A generalized framework for partial order embeddings. In *The Twelfth International Conference on Learning Representa*tions, 2024.
 - Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it?, 2023. URL https://arxiv.org/abs/2210.01936.
 - Yuhui Zhang, Yuchang Su, Yiming Liu, and Serena Yeung-Levy. Negvqa: Can vision language models understand negation?, 2025. URL https://arxiv.org/abs/2505.22946.
 - Zhanqiu Zhang, Jie Wang, Jiajun Chen, Shuiwang Ji, and Feng Wu. Cone: Cone embeddings for multi-hop reasoning over knowledge graphs. volume 34, pp. 19172–19183, 2021.

A ERROR BOUNDS FOR SOFT PROJECTORS

Lemma A.1. Let $X = U\Sigma V^{\top}$, where $U \in O(d)$, and define $\tilde{P} := X(X^{\top}X + \lambda I)^{-1}X^{\top}$. We have

$$\tilde{\boldsymbol{P}} = \boldsymbol{U}\boldsymbol{\Sigma}^2(\boldsymbol{\Sigma}^2 + \lambda \boldsymbol{I})^{-1}\boldsymbol{U}^{\top}$$
(12)

and the spectrum of $\tilde{\pmb{P}}$ is $\left\{ \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \right\}_{i=1}^d$.

Proof. Letting the SVD of X be $U\Sigma V^{\top}$,

 $\tilde{P} = U\Sigma V^{\top} (V\Sigma^{2}V^{\top} + \lambda I)^{-1}V\Sigma U^{\top}$ $= U\Sigma V^{\top}V(\Sigma^{2} + \lambda I)^{-1}V^{\top}V\Sigma U^{\top}$ $= U\Sigma^{2}(\Sigma^{2} + \lambda I)^{-1}U^{\top}.$ (13)

The spectrum of \tilde{P} is the diagonal of $\Sigma^2(\Sigma^2 + \lambda I)^{-1}$, which reads $\left\{\frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right\}_{i=1}^d$.

Proposition A.2 (Frobenius norm error). Let $X = U\Sigma V^{\top}$ be rank-r, where $U \in O(d)$ and let the orthogonal projector onto $\mathrm{Span}(X)$ be $P = UJ_rU^{\top}$. Define $\tilde{P} := X(X^{\top}X + \lambda I)^{-1}X^{\top}$. Then,

$$\|\boldsymbol{P} - \tilde{\boldsymbol{P}}\|_F \le \frac{\lambda}{\sigma_r^2 + \lambda}.\tag{14}$$

Proof. Let $J_r = \operatorname{BlockDiag}(I_r, \mathbf{0}_{d-r})$, where $r = \operatorname{rank}(\boldsymbol{X})$ and write the SVD of the orthogonal projector as $\boldsymbol{P} = \boldsymbol{U}\boldsymbol{J}_r\boldsymbol{U}^{\top}$ for $\boldsymbol{U} \in O(d)$. Using Lemma A.1, we can write $\boldsymbol{P} - \tilde{\boldsymbol{P}} = \boldsymbol{U}(\boldsymbol{J}_r - \boldsymbol{\Sigma}^2(\boldsymbol{\Sigma}^2 + \lambda \boldsymbol{I})^{-1})\boldsymbol{U}^{\top}$. The Frobenius norm is invariant to orthogonal transformations \boldsymbol{U} , hence

$$\|\boldsymbol{P} - \tilde{\boldsymbol{P}}\|_F^2 = \|\boldsymbol{J}_r - \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma}^2 + \lambda \boldsymbol{I})^{-1}\|_F^2 = \sum_{i=1}^r \left(1 - \frac{\sigma_i^2}{\sigma_i^2 + \lambda}\right)^2 \le r \left(\frac{\lambda}{\sigma_r^2 + \lambda}\right)^2.$$
 (15)

Therefore,
$$\| {m P} - \tilde{{m P}} \|_F \leq rac{\lambda \sqrt{r}}{\sigma_r^2 + \lambda}.$$

Proposition A.3 (Operator norm error). Let $X = U\Sigma V^{\top}$ be rank-r, where $U \in O(d)$, and let the orthogonal projector onto $\mathrm{Span}(X)$ be $P = UJ_rU^{\top}$. Define $\tilde{P} := X(X^{\top}X + \lambda I)^{-1}X^{\top}$. Then,

$$\|\boldsymbol{P} - \tilde{\boldsymbol{P}}\|_2 = \frac{\lambda}{\sigma_r^2 + \lambda}.\tag{16}$$

Proof. Let $J_r = \operatorname{BlockDiag}(I_r, \mathbf{0}_{d-r})$, where $r = \operatorname{rank}(\boldsymbol{X})$ and write the SVD of the orthogonal projector as $\boldsymbol{P} = \boldsymbol{U}\boldsymbol{J}_r\boldsymbol{U}^{\top}$ for $\boldsymbol{U} \in O(d)$. Using Lemma A.1, we can write $\boldsymbol{P} - \tilde{\boldsymbol{P}} = \boldsymbol{U}(\boldsymbol{J}_r - \boldsymbol{\Sigma}^2(\boldsymbol{\Sigma}^2 + \lambda \boldsymbol{I})^{-1})\boldsymbol{U}^{\top}$. The operator norm is invariant to orthogonal transformations \boldsymbol{U} , hence

$$\|\mathbf{P} - \tilde{\mathbf{P}}\|_{2} = \|\mathbf{J}_{r} - \mathbf{\Sigma}^{2} (\mathbf{\Sigma}^{2} + \lambda \mathbf{I})^{-1}\|_{2}^{2} = \max \left\{ 1 - \frac{\sigma_{i}^{2}}{\sigma_{i}^{2} + \lambda} \right\}_{i=1}^{r} = \frac{\lambda}{\sigma_{r}^{2} + \lambda}.$$
 (17)

Therefore,
$$\|P - \tilde{P}\|_2 = \frac{\lambda}{\sigma_+^2 + \lambda}$$
.

Corollary A.4 (Negation operator error). Let X, P and \tilde{P} be in the conditions of Proposition A.2. Then,

$$\|(\boldsymbol{I} - \boldsymbol{P}) - (\boldsymbol{I} - \tilde{\boldsymbol{P}})\|_2 \le \frac{\lambda}{\sigma_n^2 + \lambda}.$$
 (18)

Proof. Note that $\|(I - P) - (I - \tilde{P})\|_2 = \|P - \tilde{P}\|_2$ and apply Proposition A.3.

Proposition A.5 (Trace error). Let $X = U\Sigma V^{\top}$ be rank-r, where $U \in O(d)$, and let the orthogonal projector onto $\mathrm{Span}(X)$ be $P = UJ_rU^{\top}$. Define $\tilde{P} := X(X^{\top}X + \lambda I)^{-1}X^{\top}$. Then,

$$\left| \operatorname{Tr}(\boldsymbol{P}) - \operatorname{Tr}(\tilde{\boldsymbol{P}}) \right| \le \frac{\lambda r}{\sigma_r^2 + \lambda}$$
 (19)

Proof. First write $P - \tilde{P} = U(J_r - \Sigma^2(\Sigma^2 + \lambda I)^{-1})U^{\top}$. We have then,

$$\left| \operatorname{Tr}(\boldsymbol{P}) - \operatorname{Tr}(\tilde{\boldsymbol{P}}) \right| = \left| \operatorname{Tr}(\boldsymbol{J}_r - \boldsymbol{\Sigma}^2 (\boldsymbol{\Sigma}^2 + \lambda \boldsymbol{I})^{-1}) \right| = \sum_{i=1}^r \left(1 - \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \right) \le \frac{\lambda r}{\sigma_r^2 + \lambda}.$$
 (20)

Corollary A.6 (Subspace rank error). Letting $r := \operatorname{rank}(X)$, the relative error of estimating r via $\operatorname{Tr}(\tilde{P})$ verifies

$$\frac{\left|\text{Tr}(\tilde{\boldsymbol{P}}) - r\right|}{r} \le \frac{\lambda}{\sigma_r^2 + \lambda}.$$
 (21)

Proof. Suffices to note that $r = \text{Tr}(\mathbf{P})$ and use Proposition A.5.

Proposition A.7 (Subspace similarity error). Let X_i and X_j be rank- r_i and r_j matrices, with singular values $\{\sigma_k\}_{k=1}^d$ and $\{\eta_k\}_{k=1}^d$ (in descending order), respectively. Denote by P_i , \tilde{P}_i and P_i , \tilde{P}_i the respective orthogonal and soft projectors. Then,

$$\left| \operatorname{Tr}(\boldsymbol{P}_{i}\boldsymbol{P}_{j}) - \operatorname{Tr}(\tilde{\boldsymbol{P}}_{i}\tilde{\boldsymbol{P}}_{j}) \right| \leq \sqrt{r_{i}r_{j}} \left(\frac{\lambda}{\sigma_{r}^{2} + \lambda} + \frac{\lambda}{\eta_{r}^{2} + \lambda} \frac{\sigma_{r}^{2}}{\sigma_{r}^{2} + \lambda} \right)$$
(22)

Proof. We have

$$\left| \operatorname{Tr}(\boldsymbol{P}_{i}\boldsymbol{P}_{j}) - \operatorname{Tr}(\tilde{\boldsymbol{P}}_{i}\tilde{\boldsymbol{P}}_{j}) \right| = \left| \operatorname{Tr}((\boldsymbol{P}_{i} - \tilde{\boldsymbol{P}}_{i})\boldsymbol{P}_{j}) + \operatorname{Tr}((\boldsymbol{P}_{j} - \tilde{\boldsymbol{P}}_{j})\tilde{\boldsymbol{P}}_{i}) \right|$$

$$\leq \left| \operatorname{Tr}((\boldsymbol{P}_{i} - \tilde{\boldsymbol{P}}_{i})\boldsymbol{P}_{j}) \right| + \left| \operatorname{Tr}((\boldsymbol{P}_{j} - \tilde{\boldsymbol{P}}_{j})\tilde{\boldsymbol{P}}_{i}) \right|.$$
(23)

Apply Cauchy-Schwartz to both terms, we arrive at

$$|\operatorname{Tr}(\boldsymbol{P}_{i}\boldsymbol{P}_{j}) - \operatorname{Tr}(\tilde{\boldsymbol{P}}_{i}\tilde{\boldsymbol{P}}_{j})| \leq ||\boldsymbol{P}_{i} - \tilde{\boldsymbol{P}}_{i}||_{F}||\boldsymbol{P}_{j}||_{F} + ||\boldsymbol{P}_{j} - \tilde{\boldsymbol{P}}_{j}||_{F}||\tilde{\boldsymbol{P}}_{i}||_{F}$$
(24)

and we can replace $\sqrt{r_j} = \|P_j\|_F$, $\sqrt{r_i} = \|P_i\|_F$ and employ Proposition A.2,

$$\left| \operatorname{Tr}(\boldsymbol{P}_{i}\boldsymbol{P}_{j}) - \operatorname{Tr}(\tilde{\boldsymbol{P}}_{i}\tilde{\boldsymbol{P}}_{j}) \right| \leq \|\boldsymbol{P}_{i} - \tilde{\boldsymbol{P}}_{i}\|_{F} \sqrt{r_{j}} + \|\boldsymbol{P}_{j} - \tilde{\boldsymbol{P}}_{j}\|_{F} \|\tilde{\boldsymbol{P}}_{i}\|_{F}$$

$$\leq \sqrt{r_{i}r_{j}} \frac{\lambda}{\sigma^{2} + \lambda} + \sqrt{r_{j}} \frac{\lambda}{n^{2} + \lambda} \|\tilde{\boldsymbol{P}}_{i}\|_{F}. \tag{25}$$

Finally, note that $\|\tilde{\boldsymbol{P}}_i\|_F \leq \sqrt{r_i} \left(\frac{\sigma_r^2}{\sigma_r^2 + \lambda}\right)$

$$\left| \operatorname{Tr}(\boldsymbol{P}_{i}\boldsymbol{P}_{j}) - \operatorname{Tr}(\tilde{\boldsymbol{P}}_{i}\tilde{\boldsymbol{P}}_{j}) \right| \leq \sqrt{r_{i}r_{j}} \left(\frac{\lambda}{\sigma_{r}^{2} + \lambda} + \frac{\lambda}{\eta_{r}^{2} + \lambda} \frac{\sigma_{r}^{2}}{\sigma_{r}^{2} + \lambda} \right). \tag{26}$$

Proposition A.8 (Intersection operator error). Let X_i and X_j be rank- r_i and r_j matrices, with singular values $\{\sigma_k\}_{k=1}^d$ and $\{\eta_k\}_{k=1}^d$ (in descending order), respectively. Denote by P_i , \tilde{P}_i and P_i , \tilde{P}_j the respective orthogonal and soft projectors. Then,

$$\|\mathbf{P}_{i}\mathbf{P}_{j} - \tilde{\mathbf{P}}_{i}\tilde{\mathbf{P}}_{j}\|_{2} \leq \frac{\lambda}{\sigma_{x}^{2} + \lambda} + \frac{\lambda}{\eta_{x}^{2} + \lambda}.$$
(27)

Proof. From writing $P_iP_j-\tilde{P}_i\tilde{P}_j=(P_i-\tilde{P}_i)P_j+(P_j-\tilde{P}_j)\tilde{P}_i$ and applying the triangle inequality

$$\|P_{i}P_{j} - \tilde{P}_{i}\tilde{P}_{j}\|_{2} = \|(P_{i} - \tilde{P}_{i})P_{j} + (P_{j} - \tilde{P}_{j})\tilde{P}_{i}\|_{2}$$

$$\leq \|P_{i} - \tilde{P}_{i}\|_{2}\|P_{j}\|_{2} + \|P_{j} - \tilde{P}_{j}\|_{2}\|\tilde{P}_{i}\|_{2}.$$
(28)

Noting that $\|P_i\|_2 \le 1$ and $\|\tilde{P}_i\|_2 \le 1$ and using Proposition A.3, we have

$$\|\mathbf{P}_i \mathbf{P}_j - \tilde{\mathbf{P}}_i \tilde{\mathbf{P}}_j\|_2 \le \frac{\lambda}{\sigma_z^2 + \lambda} + \frac{\lambda}{n_z^2 + \lambda}.$$
 (29)

B GRADIENTS OF SOFT PROJECTION MATRICES

To understand how gradient-based training inherently shapes subspaces, we analyze the gradient flow of the subspace intersection. This reveals how projection operators evolve by incorporating missing dimensions from positive samples and repelling those aligned with negative ones.

The gradient of $\text{Tr}(\tilde{P}_i\tilde{P}_j)$ with respect to X_i can be derived from the identity

$$\nabla_{\boldsymbol{X}} \operatorname{Tr} \left((\boldsymbol{A} + \boldsymbol{X}^{\top} \boldsymbol{X} \boldsymbol{X})^{-1} (\boldsymbol{X}^{\top} \boldsymbol{B} \boldsymbol{X}) \right) =$$

$$-2C\boldsymbol{X} (\boldsymbol{A} + \boldsymbol{X}^{\top} \boldsymbol{C} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} \boldsymbol{B} \boldsymbol{X} (\boldsymbol{A} + \boldsymbol{X}^{\top} \boldsymbol{C} \boldsymbol{X})^{-1} + 2\boldsymbol{B} \boldsymbol{X} (\boldsymbol{A} + \boldsymbol{X}^{\top} \boldsymbol{C} \boldsymbol{X})^{-1}.$$
 (30)

We have then

$$\nabla_{\boldsymbol{X}_{i}} \operatorname{Tr} \left(\tilde{\boldsymbol{P}}_{i} \tilde{\boldsymbol{P}}_{j} \right) = 2(\boldsymbol{I} - \tilde{\boldsymbol{P}}_{i}) \tilde{\boldsymbol{P}}_{j} \boldsymbol{X}_{i} (\boldsymbol{X}_{i}^{\top} \boldsymbol{X}_{i} + \lambda \boldsymbol{I})^{-1}$$

$$\propto \underbrace{\tilde{\boldsymbol{P}}_{i}^{\perp} \tilde{\boldsymbol{P}}_{j}}_{\text{New information}} \underbrace{\boldsymbol{X}_{i} (\boldsymbol{X}_{i}^{\top} \boldsymbol{X}_{i} + \lambda \boldsymbol{I})^{-1}}_{\text{Spectral scaling}}.$$
(31)

The spectral scaling factor $\boldsymbol{X}_i(\boldsymbol{X}_i^{\top}\boldsymbol{X}_i+\lambda\boldsymbol{I})^{-1}$ acts as low-pass filter on \boldsymbol{X}_i . If we write the SVD of \boldsymbol{X}_i as $\boldsymbol{X}_i=\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\top}$, then $\boldsymbol{X}_i(\boldsymbol{X}_i^{\top}\boldsymbol{X}_i+\lambda\boldsymbol{I})^{-1}=\boldsymbol{U}\boldsymbol{\Sigma}(\boldsymbol{\Sigma}^2+\lambda\boldsymbol{I})^{-1}\boldsymbol{V}^{\top}$. As a result, high-energy directions (associated with large singular values) are attenuated, while low-energy directions are amplified. This ensures that updates to \boldsymbol{X}_i preserve dominant, well-supported directions while adapting underrepresented ones.

The component $\tilde{P}_i^\perp \tilde{P}_j$, where $\tilde{P}_i^\perp = I - \tilde{P}_i$, indicates that gradient flow occurs only along directions present in subspace j but orthogonal to subspace i, formally, in $\mathrm{range}(\tilde{P}_j) \cap \mathrm{null}(\tilde{P}_i)$. Thus, the learning signal drives X_i to incorporate directions it lacks but that are represented by X_j , encouraging alignment without redundancy. If subspace j is already contained within subspace i i.e., $\tilde{P}_j \leq \tilde{P}_i$, the gradient vanishes since $\tilde{P}_j \tilde{P}_i = \tilde{P}_j$ implies $(I_d - \tilde{P}_i) \tilde{P}_j = 0$. This update mechanism shares similarities with Oja's rule in online PCA, promoting efficient subspace adaptation.

Conversely, negative pairs induce repulsive gradients, driving X_i to remove directions aligned with X_j and thus promoting subspace separation. Consequently, the effective dimensionality of subspace i naturally adapts to encompass the union of all its relevant positive neighbors i.e.,

$$\operatorname{rank}(\tilde{\mathbf{P}}_i) \ge \dim \operatorname{span}\left(\bigcup_{j \in \operatorname{Pos}(i)} \operatorname{range}(\tilde{\mathbf{P}}_j)\right). \tag{32}$$

In other words, examples with more diverse positive neighborhoods require richer subspaces, while simpler ones can be encoded more compactly.

C WORDNET EXPERIMENTS

WORDNET's noun hierarchy has 82,115 nodes and 75,850 edges. The verb hierarchy is smaller, featuring 13,767 nodes and 13,239 edges. Their transitive closures are significantly denser, with 663,508 (noun) and 35,079 (verb) edges. All WORDNET experiments were conducted on a RTX8000 GPU with 49GB of memory.

C.1 RECONSTRUCTION

Experimental Details. We parameterize each node's subspace with a matrix $X_i \in \mathbb{R}^{128 \times 128}$, initialized with entries from a zero-mean Gaussian distribution with standard deviation 0.0001. The regularizer was set $\lambda = 0.2$. For each training edge (u,v), we sample 19 nodes $v' \neq u$ such that neither (u,v') nor (v',u) are in the train split and optimized InfoNCE, applying the the subspace similarity $\mathrm{Tr}(\tilde{P}_i\tilde{P}_j)$ from Eq. (3) to soft projectors. We used Adam Kingma & Ba (2017), with a batch-size of 128 and learning rate of 0.0005. During evaluation, we compute the similarity $\mathrm{Tr}(\tilde{P}_u\tilde{P}_v)$ of each edge (u,v) in the full transitive closure $\mathrm{TC}(\mathcal{G})$ and rank it among the those of all node pairs that are not connected in the transitive closure $\{\mathrm{Tr}(\tilde{P}_u\tilde{P}_{v'}): (u,v') \notin \mathrm{TC}(\mathcal{G})\}$.

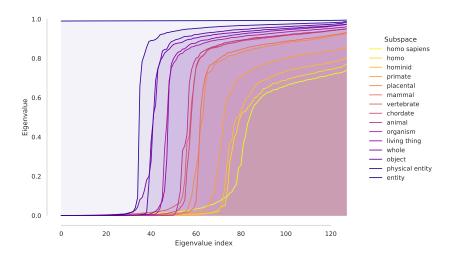


Figure 7: Sorted eigenvalues of the soft projection operators \vec{P} for nouns in the hypernymy chain homo sapiens \rightarrow entity. As we move from specific to general concepts, the subspace's effective rank gradually expands. This illustrates how our soft projectors naturally capture concept specificity: specific nouns like homo sapiens activate fewer dimensions (eigenvalues near zero), while broader concepts like entity activate more dimensions (eigenvalues near one).

Visualization of the Spectrum of WordNet Nouns. In Fig. 7, we plot the sorted eigenvalues of our soft projector representations (\tilde{P}) for WORDNET nouns, traversing a hypernymy chain from homo sapiens to entity. This plot illustrates two key properties:

- Smooth Eigenvalue Distribution: Unlike the binary (0 or 1) eigenvalues of orthogonal projection matrices, the eigenvalues of \tilde{P} are smooth within [0,1]. This smoothness is crucial for our learnable, soft subspace representations.
- Effective Rank Justification: The plot directly justifies our use of $\operatorname{Tr}(\tilde{P})$ as a measure of the *effective rank* of a concept's subspace. For orthogonal projection operators, the trace (sum of eigenvalues) precisely equals the subspace's rank due to their binary eigenvalues. Here, while eigenvalues are not binary, the plot clearly shows the distribution of activated dimensions for each concept. For instance, the broad concept *entity* utilizes all 128 dimensions, with all eigenvalues near one. In contrast, *homo sapiens* activates fewer dimensions, with most eigenvalues approaching zero.

C.2 LINK PREDICTION

Experimental Details. For link prediction, every node is initialized as a random matrix $X_i^{d\times n}$, with entries sampled from a zero-mean Gaussian distribution ($\sigma=0.0001$). In our experiments we considered d=n=64 as well as d=n=128. The soft projector regularizer was set to $\lambda=0.2$. We optimized the margin loss from Eq. (8) with $\gamma_+=0.9$ and $\gamma_-=0.5$ for 0% of non-basic edges, and $\gamma_+=0.8$, $\gamma_-=0.1$ for the remaining percentages. To compute this loss, we used 10 negatives per each observed positive edge (u,v). Negatives were generated by sampling 5 corrupted-tail (u,v') and 5 corrupted-head (u',v) examples per positive edge, with corrupted nodes sampled from the entire set of nodes. We employed Adam (Kingma & Ba, 2017) with a constant learning rate of 0.0005 and a batch-size of 128 to perform the optimization.

C.3 GRADED LEXICAL ENTAILMENT

Experimental Details. For our HYPERLEX experiment, we use the noun subset (2,163 pairs), which provides human-annotated scores (0-10) for word pairs (u,v), quantifying the degree to which u is a type of v. We quantify entailment using the NIS from Eq. (4), with word sense disambiguation performed as in Athiwaratkun & Wilson (2018), by selecting the WORDNET synset pair with maximal subspace similarity $\text{Tr}(\tilde{P}_i, \tilde{P}_j)$.

D NLI EXPERIMENTS

Experimental Details. All experiments utilized a maximum sequence length of 35. We trained all-MiniLM-L6-v2 with a batch size of 1024, and all-mpnet-base-v2 with a batch size of 2048. Optimization was performed using Adam (Kingma & Ba, 2017), employing a learning rate of 0.0001 and no weight decay. An exponential learning-rate scheduler with a gamma of 0.9 was used. For the MLP-based baselines, premise and hypothesis embeddings were first computed by mean pooling the transformer's output hidden state before being passed to the MLP classification head. The MLP classification head consisted of 3 layers, featuring LeakyReLU activations and matching the hidden dimension of its corresponding transformer. A label smoothing of 0.1 was consistently applied across all training runs. The Beta priors of our model were initialized as ($\alpha_C = 1$, $\beta_C = 6$) and ($\alpha_E = 6$, $\beta_E = 1$) and were optimized during training. All experiments were conducted on a RTX8000 GPU with 49GB of memory.

E COMPOSITE QUERIES

In this section, we present in more detail how to construct logical queries with subspace operations. We use our SNLI-fine-tuned mpnet-base-v2 + SPH ($\rm SE^{128}$) model to embed the 155,070 Flickr30k captions (as candidate sentences) and encode individual phrases for logical query construction. Composite query subspaces are then formed by combining these individual embeddings.

For instance, to represent the concept "a dog running", we obtain its subspace embedding $\mathcal{S}_{\text{a dog running}}$, which is represented, in practice, by the output of the SPH, a soft projector $\tilde{\boldsymbol{P}}_{\text{a dog running}} \in \mathbb{R}^{128 \times 128}$. Similarly for "on the beach", we obtain $\tilde{\boldsymbol{P}}_{\text{on the beach}} \in \mathbb{R}^{128 \times 128}$. These subspaces are then combined using linear-algebraic operations to formulate composite query subspaces which approximate the true logical connectives described in the main text.

Conjunction. We approximate the projection operator of subspace intersection $S_i \cap S_j$ as $\tilde{P}_i \tilde{P}_j$. Recall that $P_i P_i$ is the projector onto the intersection if and only if P_i and P_j commute *i.e.*, $P_i P_j = P_j P_i$. While this is not enforced in our training pipeline, we observed good empirical results from approximating the intersection operator by $\tilde{P}_i \tilde{P}_j$.

Negation. We approximate the projection operator onto the subspace complement S_i^{\perp} as $I - \tilde{P}_i$. The approximation error in this case only comes from the regularization λ .

Disjunction. We approximate the projection operator onto the subspace sum $S_i + S_j$ by explicitly building the linear sum from the basis of each soft projector. Letting the SVD of the projectors be $\tilde{P}_i = U_i \Sigma_i V_i^{\top}$ and $\tilde{P}_j = U_j \Sigma_j V_j^{\top}$, we approximate the soft projector for the subspace sum as $X(X^{\top}X + \lambda I)^{-1}X^{\top}$, where $X = \begin{bmatrix} U_i \Sigma_i^{\frac{1}{2}} & U_j \Sigma_j^{\frac{1}{2}} \end{bmatrix}$. Here, $\lambda = 0.05$ (consistent with the regularization used during training).

F EFFICIENCY EXPERIMENTS

Retrieval Efficiency. We benchmarked top-10 retrieval latency on the 155,070 captions from the Flickr30k dataset, using batches of 128 queries. We compared our subspace embeddings (SE¹²⁸) against a 10-dimensional Poincaré hyperbolic baseline (\mathcal{P}^{10}). Because hyperbolic distance is non-Euclidean, we applied brute-force search over the entire database, ranking by the negative hyperbolic distance. In contrast, our NIS score can be formulated as a maximum inner product search problem between query and caption vectors:

$$NIS(\tilde{\boldsymbol{P}}_{caption} \mid \tilde{\boldsymbol{P}}_{query}) = \left(\frac{\text{vec}(\tilde{\boldsymbol{P}}_{caption})}{\text{Tr}(\tilde{\boldsymbol{P}}_{caption})}\right)^{\top} \text{vec}(\tilde{\boldsymbol{P}}_{query}). \tag{33}$$

This formulation allows us to use fast approximate search libraries. We indexed the normalized caption vectors using a CPU index from the FAISS library Douze et al. (2025), specifically an inverted file index with Product Quantization (IndexIVFPQ). The index was trained on 50,000 vectors. We

Table 7: **GPU encoding time per query (ms)**, averaged over Flickr30k's captions. The overhead of the SPH module is consistently small.

	Batch size				
Model	1	4	16	64	128
mpnet-base-v2	5.95	1.74	0.69	0.59	0.56
mpnet-base-v2 + SPH (SE^{128})	6.80	2.13	0.86	0.73	0.68

used 64 subquantizers for PQ with 8 bits per subquantizer, and set the search-time parameter to $n_{\rm probe}=32.$

Encoding Time. We also measured the overhead introduced by our SPH module when encoding Flickr30k captions on a RTX8000 GPU with 49GB of memory. To isolate the computation cost of the forward pass, tokenization (max-size of 35) and data transfers were computed beforehand. Table 7 shows the average encoding time per query (forward-pass) for different batch sizes, demonstrating that the additional computational cost is modest, especially with larger batches.

G LARGE LANGUAGE MODELS

The authors are solely responsible for the research ideas, experimental design, and analysis presented in this work. Large language model (LLMs) was used for editorial assistance to enhance the paper's clarity and readability, with its contributions limited to grammar, sentence structure, and flow.

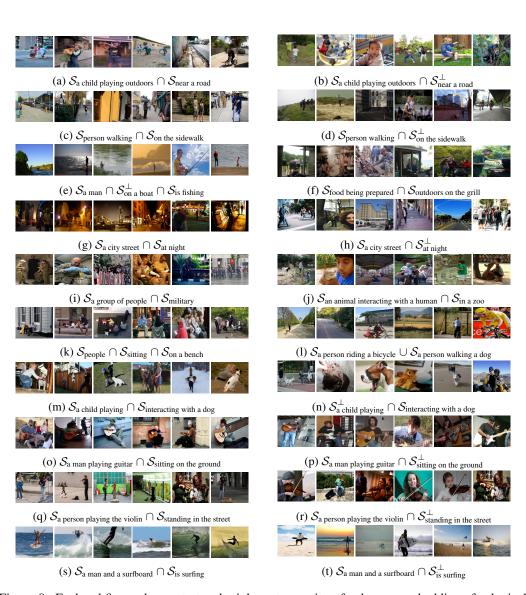


Figure 8: Each subfigure demonstrates the inherent capacity of subspace embeddings for logical composition. Queries are formed by applying subspace operations intersection (\cap), linear sum (+) and orthogonal complement (\perp) to the subspace embeddings of phrases or sentences, embedded by our SNLI-fine-tuned mpnet-base-v2 + SPH (SE¹²⁸) model. For each composite query, we retrieve the top Flickr30k images whose captions have the highest NIS with the composite query subspace.