Geometry of Categorical & Hierarchical Concepts in LLMs

Kiho Park, Yo Joong Choe, Yibo Jiang, Victor Veitch

University of Chicago

Problems

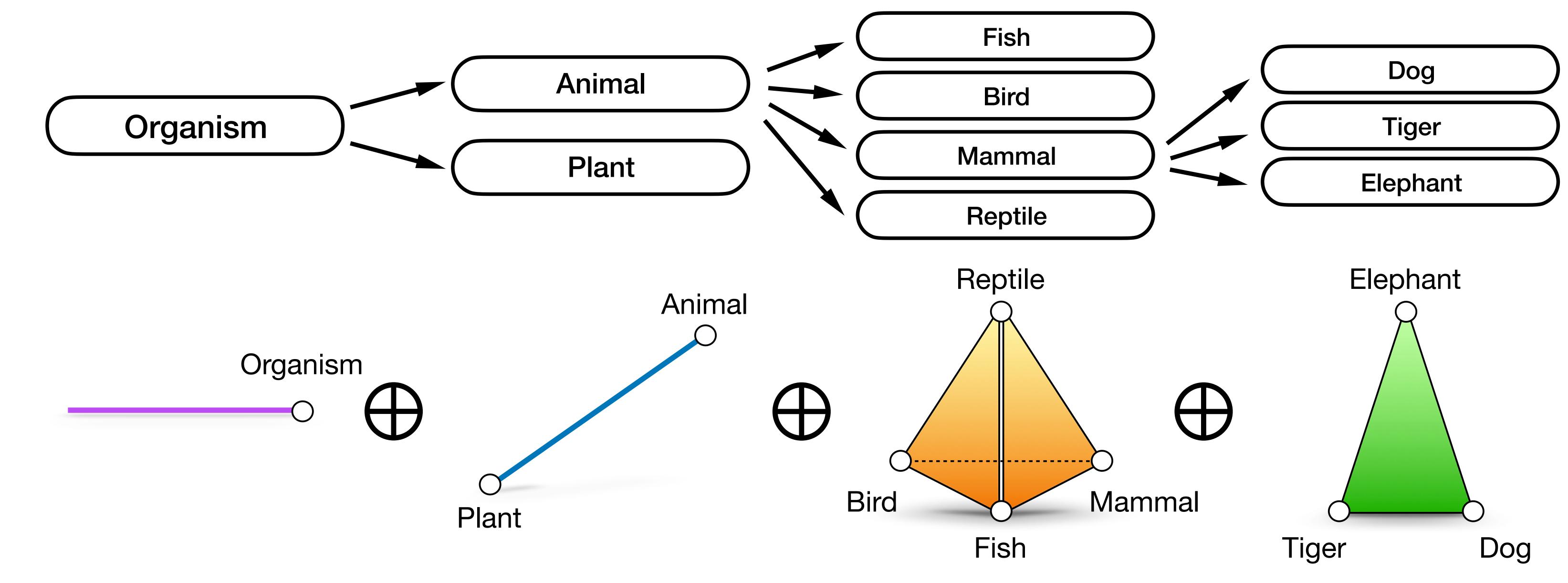
How are categorical concepts represented?

How are hierarchical relations between concepts represented?

Goal: extend the linear representation hypothesis to answer these problems

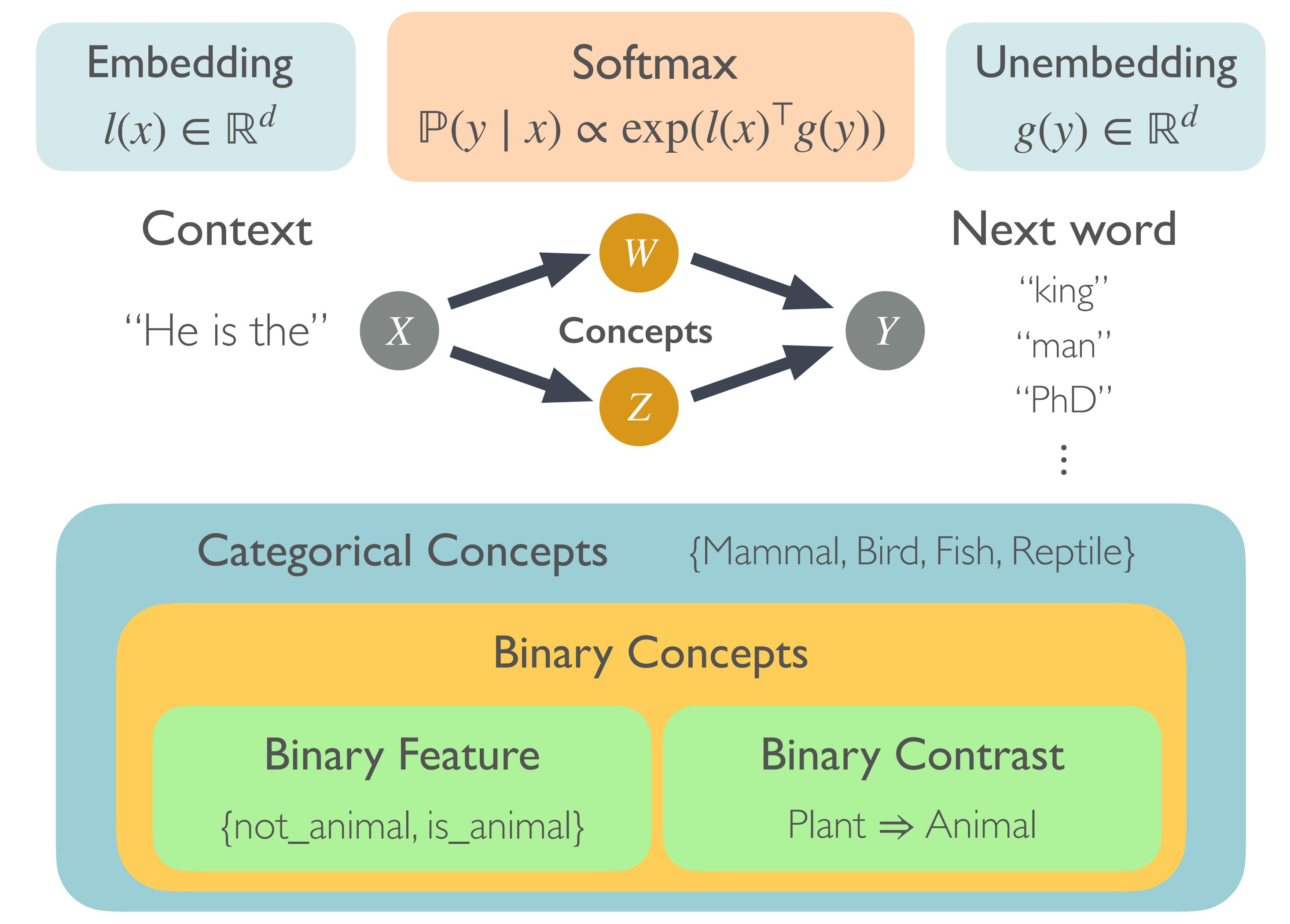
Concepts in LLMs

LLMs generate the next word using the softmax distribution. We use the unified space induced by the causal inner product.



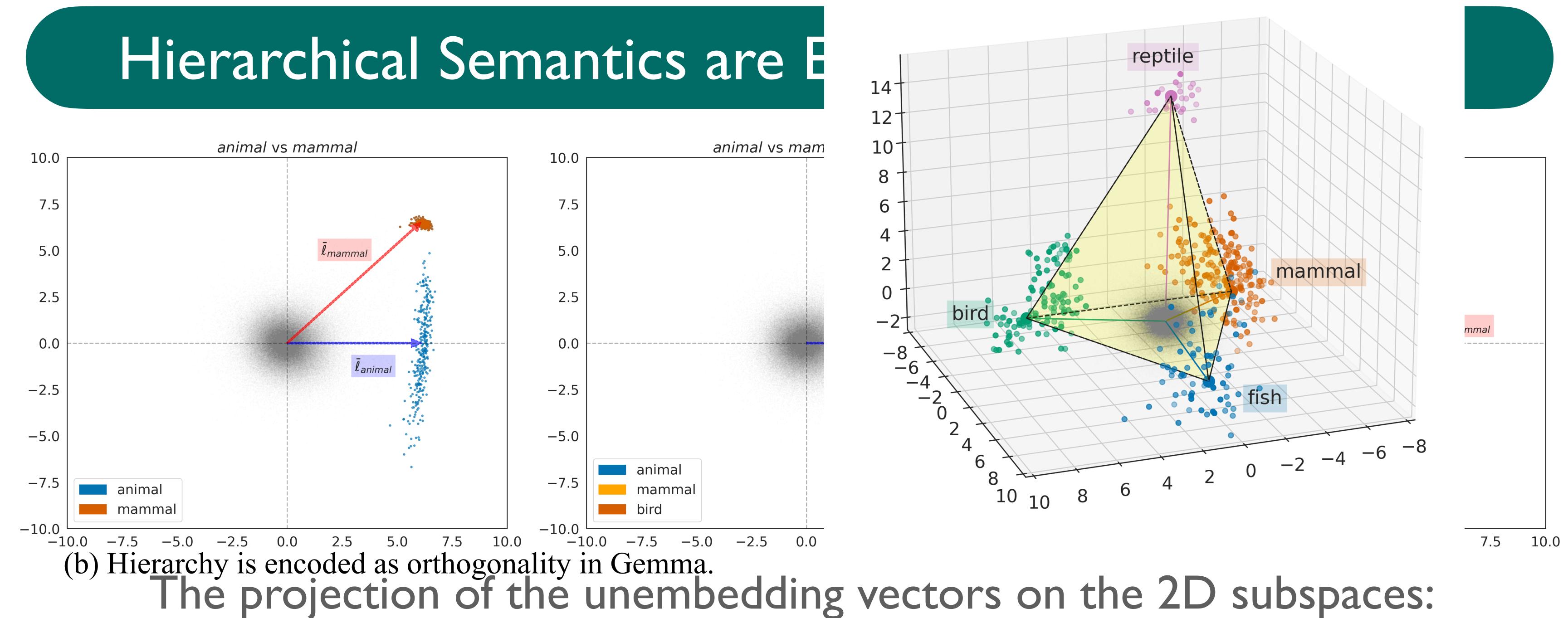
Results

Pictorial depiction of the representation of hierarchically related concepts



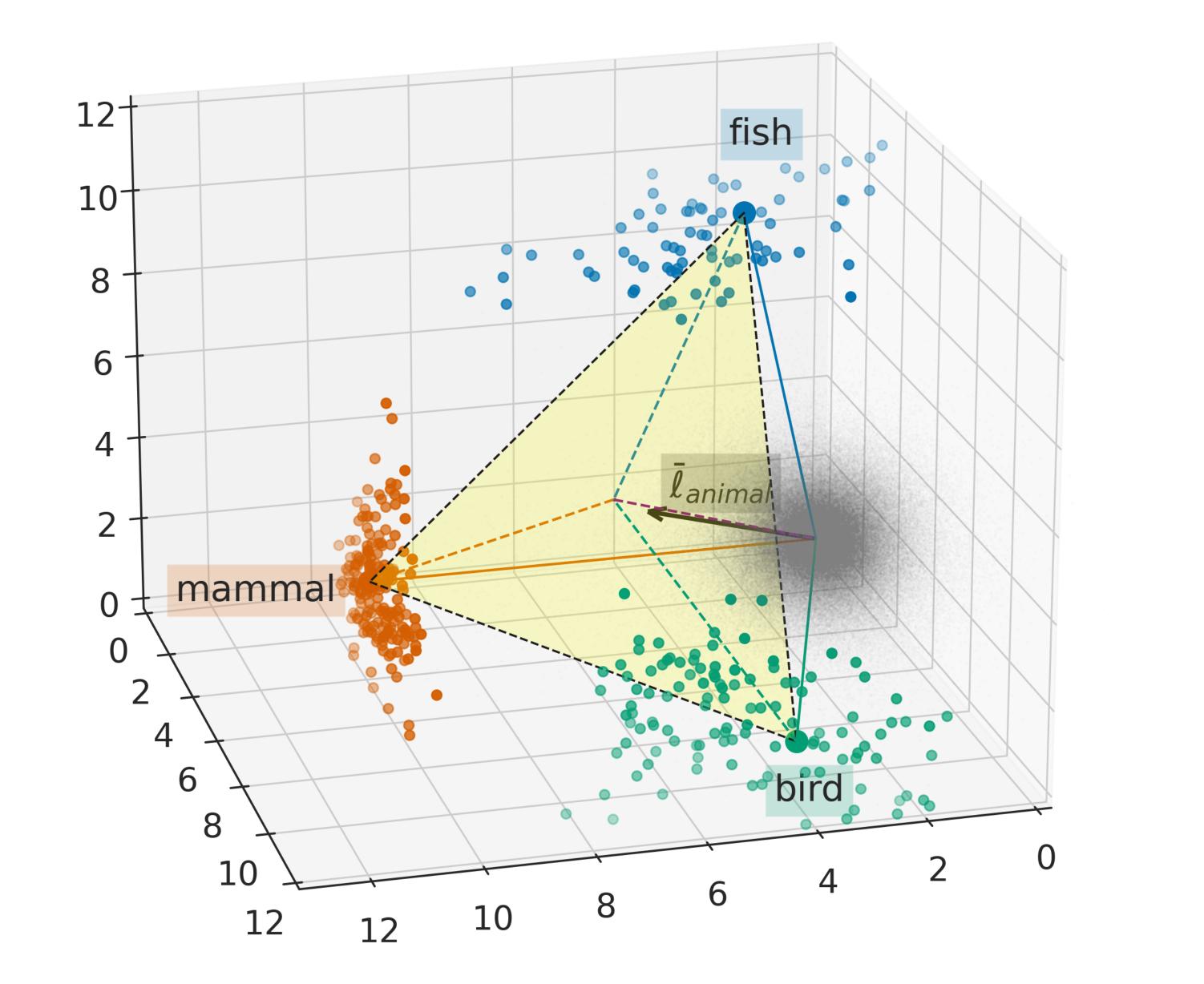
Hierarchical Structure: z is subordinate to w (denoted by $z \prec w$) if $\mathscr{Y}(z) \subseteq \mathscr{Y}(w)$. We say that Z is subordinate to W if there exists a value w_Z of W such that each value z_i of Z is subordinate to w_Z .

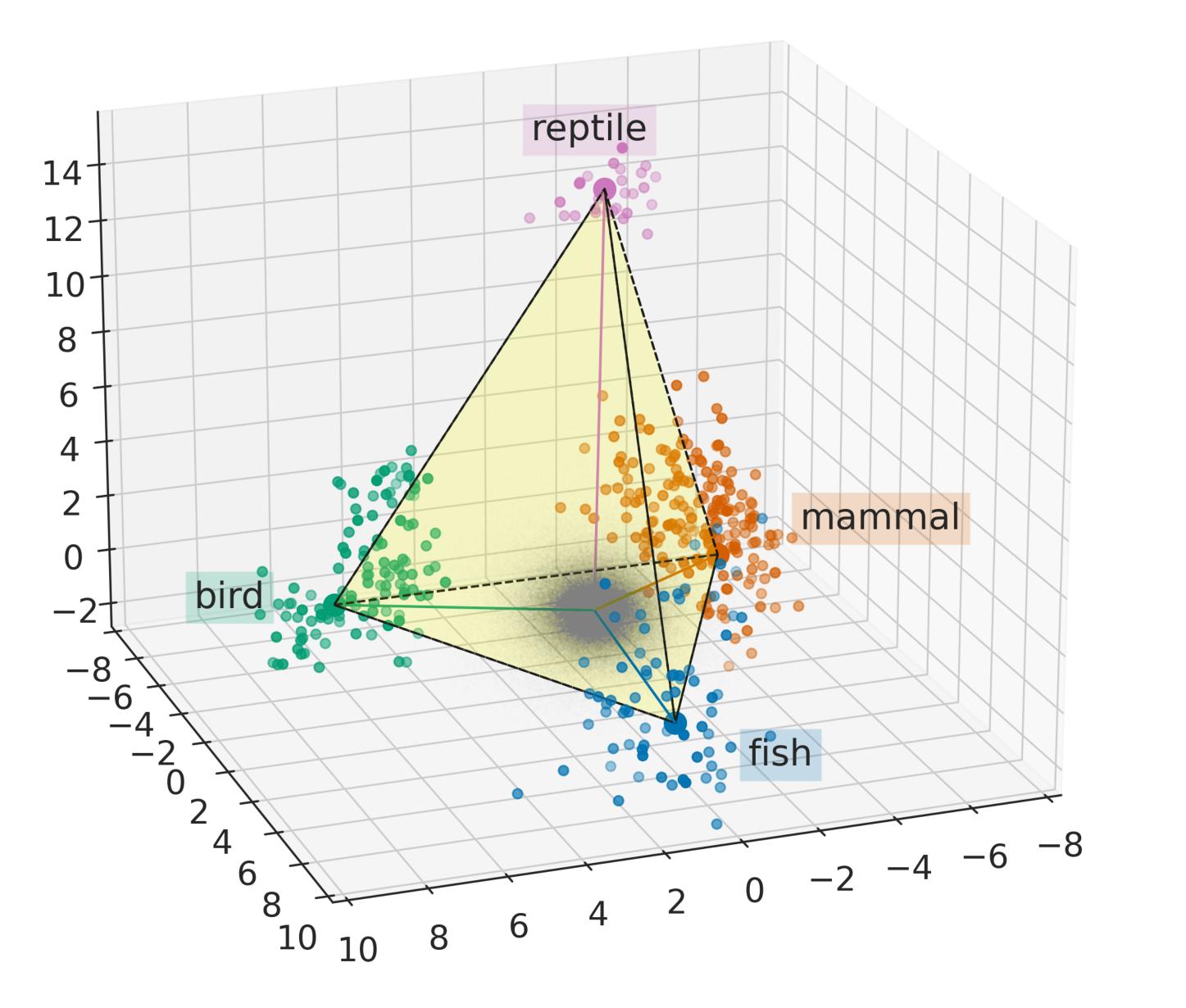
Linear Representation of Binary Concept



span{ \bar{l}_{animal} , \bar{l}_{mammal} } (left, Thm 6 (b)), span{ \bar{l}_{animal} , \bar{l}_{bird} – \bar{l}_{mammal} } (middle, Thm 6 (c)), span{ \bar{l}_{animal} – \bar{l}_{plant} , \bar{l}_{bird} – \bar{l}_{mammal} } (right, Thm 6 (d))

Categorical Concepts are Represented as Simplices





A vector \bar{l}_W is a linear representation of a binary concept W if $\mathbb{P}(W = 1 \mid l + \alpha \bar{l}_W) > \mathbb{P}(W = 1 \mid l)$ $\mathbb{P}(Z \mid l + \alpha \bar{l}_W) = \mathbb{P}(Z \mid l)$

for all contexts $l, \alpha > 0, Z$ subordinate to or causally separable with W.

Vector Representation of Binary Feature (Theorem 4)

If there exists a linear representation of a binary feature W for a value w, then, with a canonical origin, there exists a constant $b_w > 0$ such that

$$I_{W}^{\top}g(y) = \begin{cases} b_{w} & \text{if } y \in \mathscr{Y}(w) \\ 0 & \text{if } y \notin \mathscr{Y}(w) \end{cases}$$
flower
tree
dog
not_animal
is_animal

Now the linear representation has a magnitude, and we define a vector representation \bar{l}_w of a binary feature for a value w where $\|\bar{l}_w\|_2 = b_w$.

Hierarchical Orthogonality (Theorem 6)

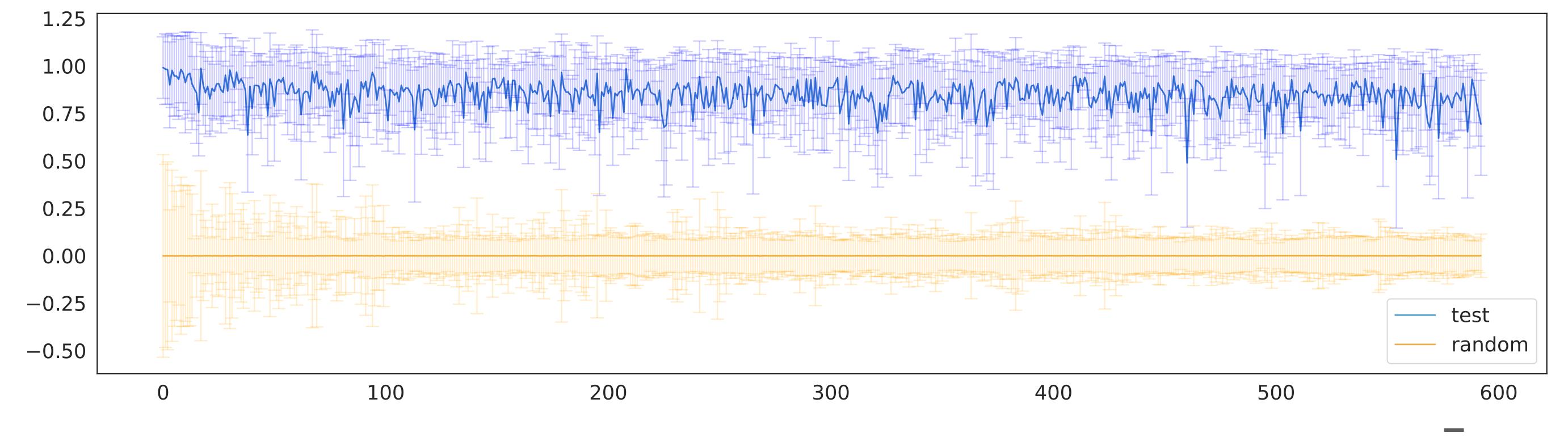
(a)
$$\bar{l}_{w_1} - \bar{l}_{w_0}$$
 is a linear representation $\bar{l}_{w_0 \Rightarrow w}$

The projection of the unembedding vectors on the 3D subspaces: span{ $\bar{l}_{mammal}, \bar{l}_{bird}, \bar{l}_{fish}$ } (left, \bar{l}_{animal} is orthogonal to the simplex!), span{ $\bar{l}_{bird} - \bar{l}_{mammal}, \bar{l}_{fish} - \bar{l}_{mammal}, \bar{l}_{reptile} - \bar{l}_{mammal}$ } (right)

Further Experiments and Details

Using Gemma-2B model, we estimate the vector representation \bar{l}_w of a binary feature for a value w using the collection of the tokens that have w:

$$\bar{l}_{w} = (\tilde{g}_{w}^{\mathsf{T}}\mathbb{E}(g_{w}))\tilde{g}_{w}, \text{ with } \tilde{g}_{w} = \frac{\operatorname{Cov}(g_{w})^{\dagger}\mathbb{E}(g_{w})}{\|\operatorname{Cov}(g_{w})^{\dagger}\mathbb{E}(g_{w})\|_{2}}$$



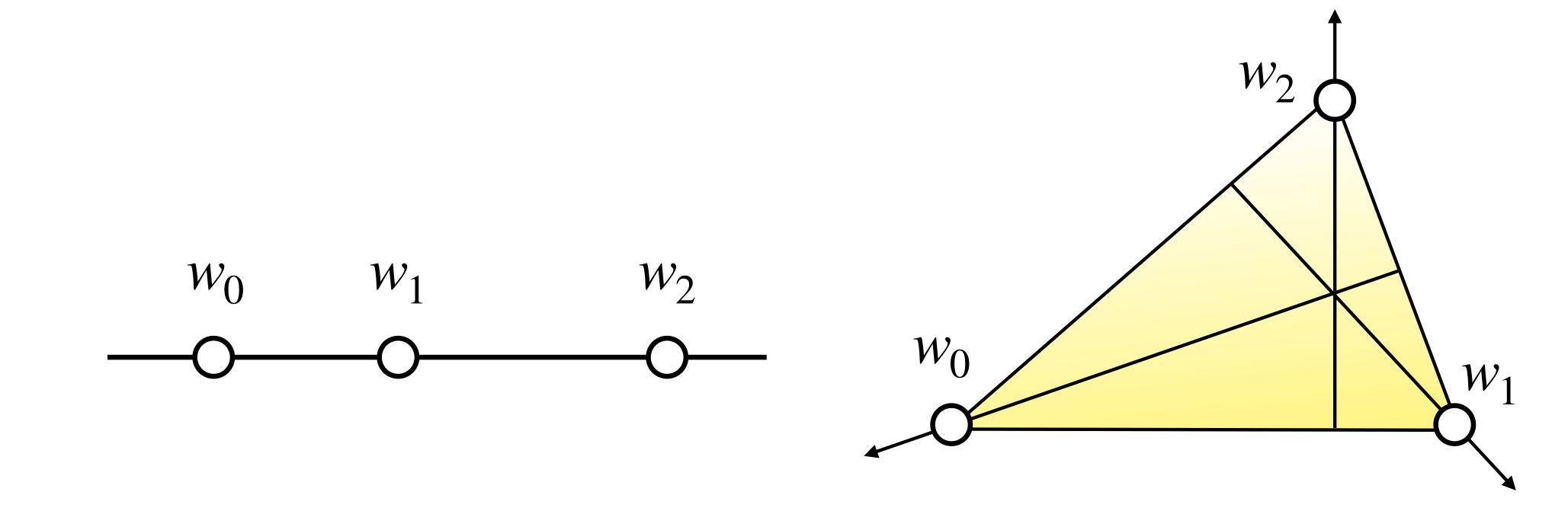
The projection of test and random words onto the estimated \overline{l}_{w} shows

(b) $\bar{l}_{w} \perp \bar{l}_{z} - \bar{l}_{w}$ for $z \prec w$ (c) $\bar{l}_{w} \perp \bar{l}_{z_{1}} - \bar{l}_{z_{0}}$ for $Z \in_{R} \{z_{0}, z_{1}\} \prec W \in_{R} \{\text{not_w, is_w}\}$ (d) $\bar{l}_{w_{1}} - \bar{l}_{w_{0}} \perp \bar{l}_{z_{1}} - \bar{l}_{z_{0}}$ for $Z \in_{R} \{z_{0}, z_{1}\} \prec W \in_{R} \{w_{0}, w_{1}\}$ (e) $\bar{l}_{w_{1}} - \bar{l}_{w_{0}} \perp \bar{l}_{w_{2}} - \bar{l}_{w_{1}}$ for $w_{2} \prec w_{1} \prec w_{0}$

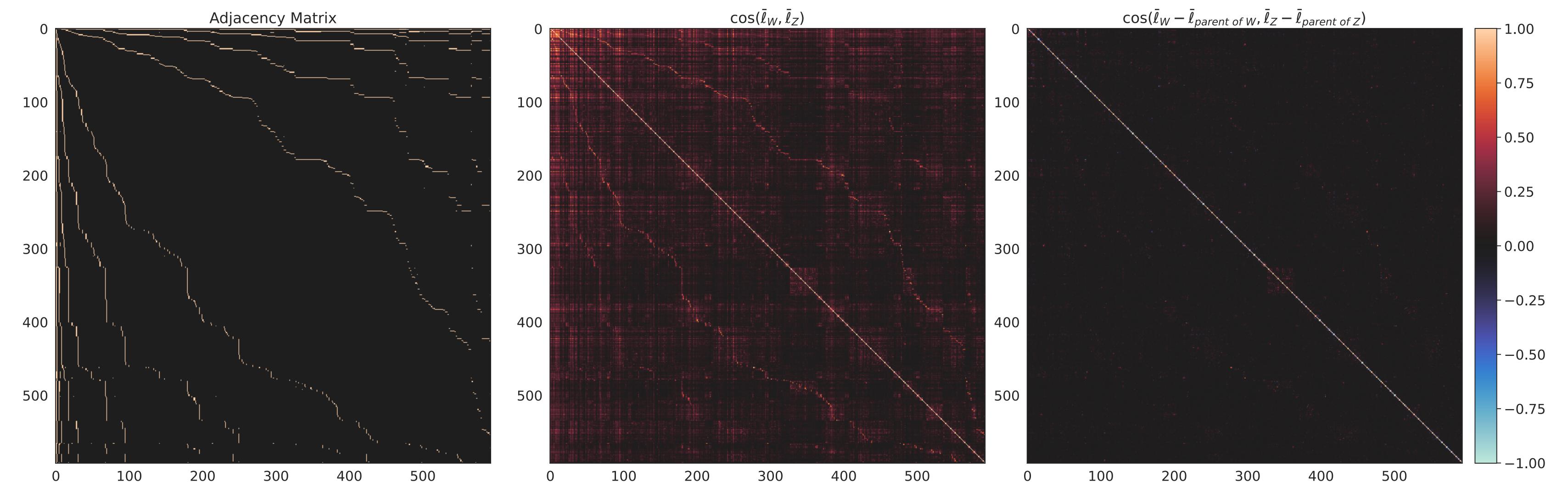
Categorical Concepts as Simplices (Theorem 8)

The polytope representation of a categorical concept Z is the convex hull of the vector representations for the elements of the concept.

For every joint distribution $Q(w_0, ..., w_{k-1})$, if there exists some l_i such that $\mathbb{P}(W = w_i \mid l_i) = Q(W = w_i)$ for every *i*, the vector representations $\overline{l}_{w_0}, ..., \overline{l}_{w_{k-1}}$ form a (k - 1)-simplex in the representation space. In this case, we take the simplex to be the representation of the categorical concept $W = \{w_0, ..., w_{k-1}\}$.

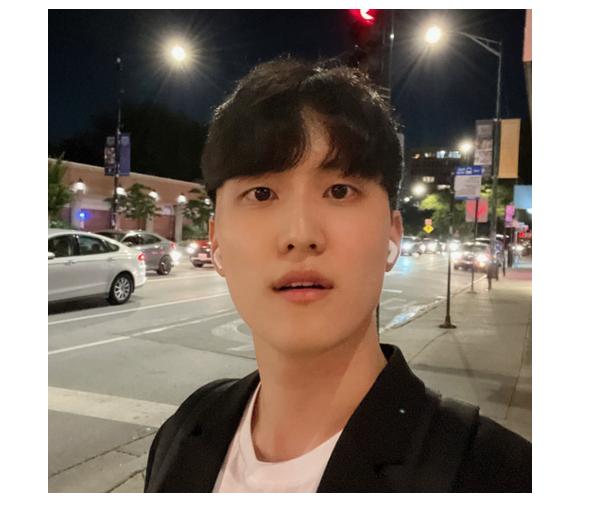


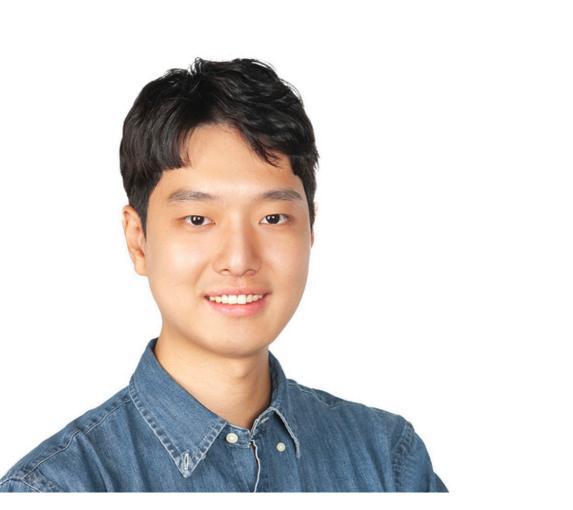
that the vector representations for binary features in WordNet exist.



The adjacency matrix of WordNet hierarchy (left) is clearly visible in the middle heatmap. By contrast, consistent with Thm 6 (e), the child-parent and parent-grandparent vectors are orthogonal in the right heatmap.

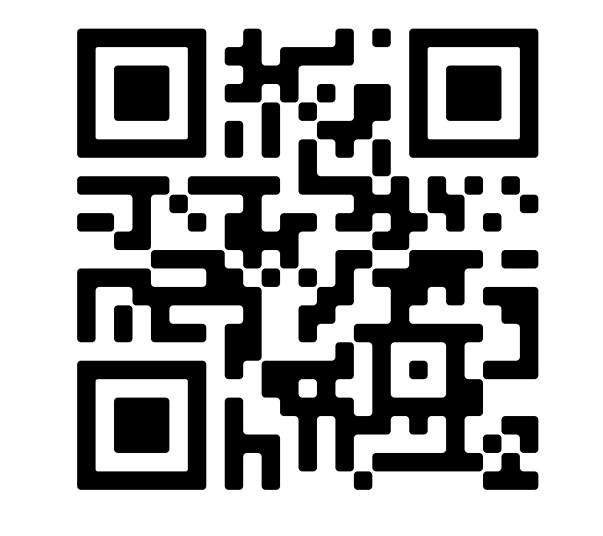
Find us online! Email: parkkiho@uchicago.edu











Kiho Park Yo

Yo Joong Choe

Yibo Jiang

Victor Veitch arXiv:2406.01506