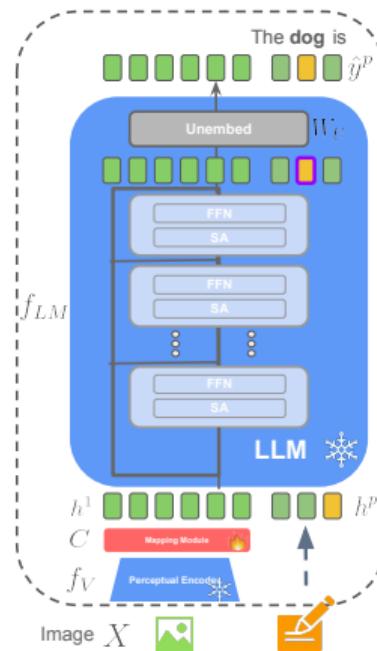
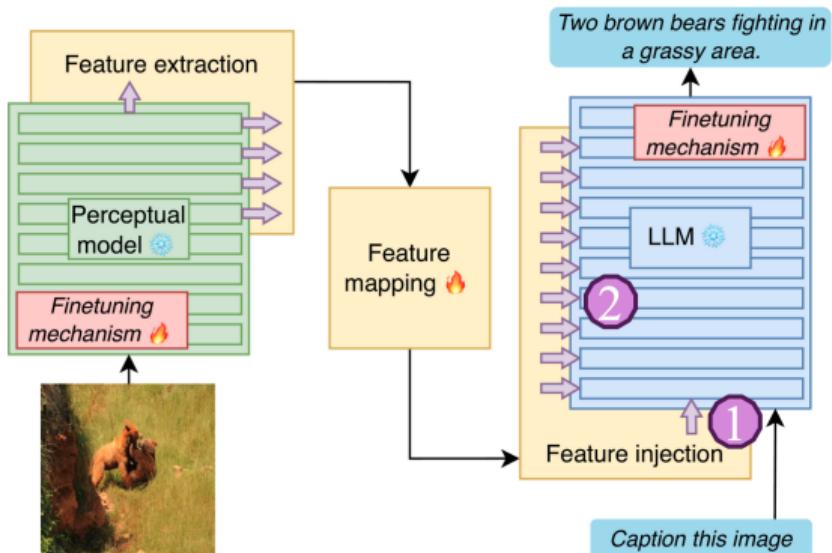


Explaining/Monitoring LMMs



Monitoring LLMs

- LMM $f = f_{LM}(c(f_v))$
long model \ connector \ visual encoder
- Understand f internal representation of a given token "dog" using CoX-LLM
 - one example $\rightarrow B$ samples $\rightarrow K$ clusters U_K
in latent space
 - U_K characterized both in terms of text and of images

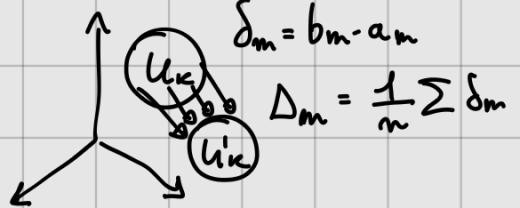
Fine-tuned LLMs

After fine-tuning concepts change

- \hookrightarrow refined
- \hookrightarrow emerged
- \hookrightarrow diminished

\hookrightarrow stronger for new concepts,
inefficient for some previously learnt concept (forgotten)

\hookrightarrow compute concepts shift



Steering

Linear modification to align the output in latent space

\hookrightarrow refusal,
truthful,...

P2S \longrightarrow L2S

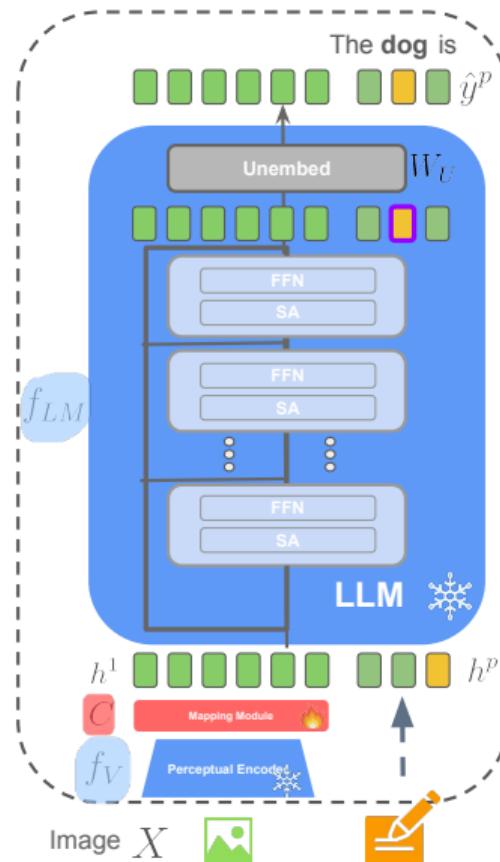
fixed steering
direction

steering direction
input dependent

Explaining/Monitoring LMMs

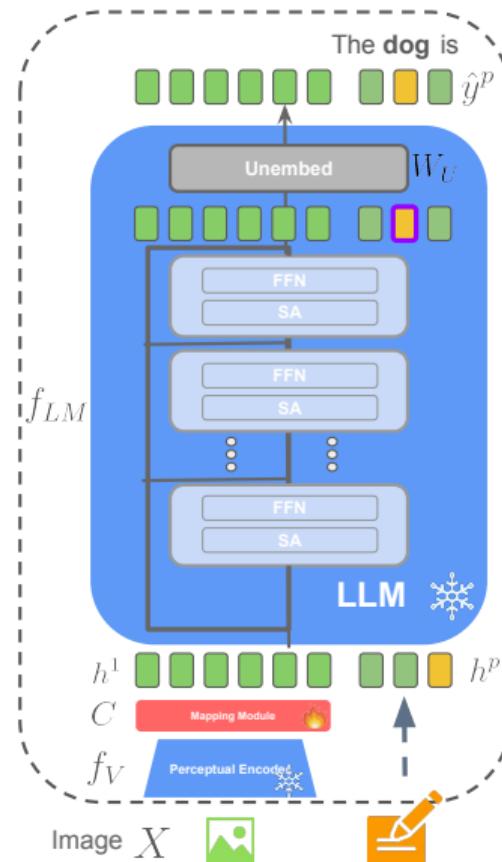
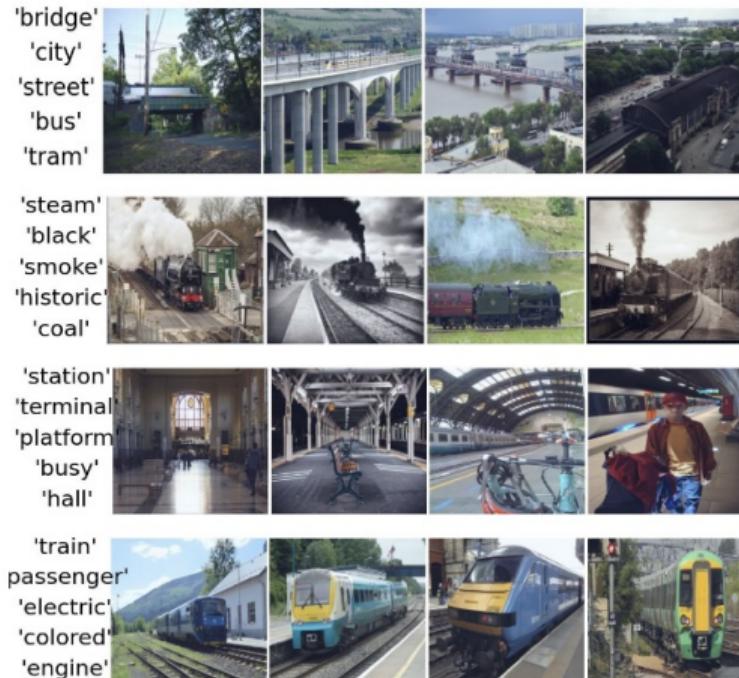
- ▶ Pretrained LMM $f = \text{Visual encoder } (f_V) + \text{Connector } (C) + \text{Language model } (f_{LM})$
- ▶ Captioning dataset $\mathcal{S} = \{(X_i, y_i)\}_{i=1}^N$.
Images $X_i \in \mathcal{X}$ and captions $y_i \subset \mathcal{Y}$
- ▶ A token of interest $t \in \mathcal{Y}$ (Eg. 'Dog', 'Cat' etc.)
- ▶ **Analysis:** Understand internal representations of f about t in terms of high-level concepts

Concept based eXplainability framework
for LMMs

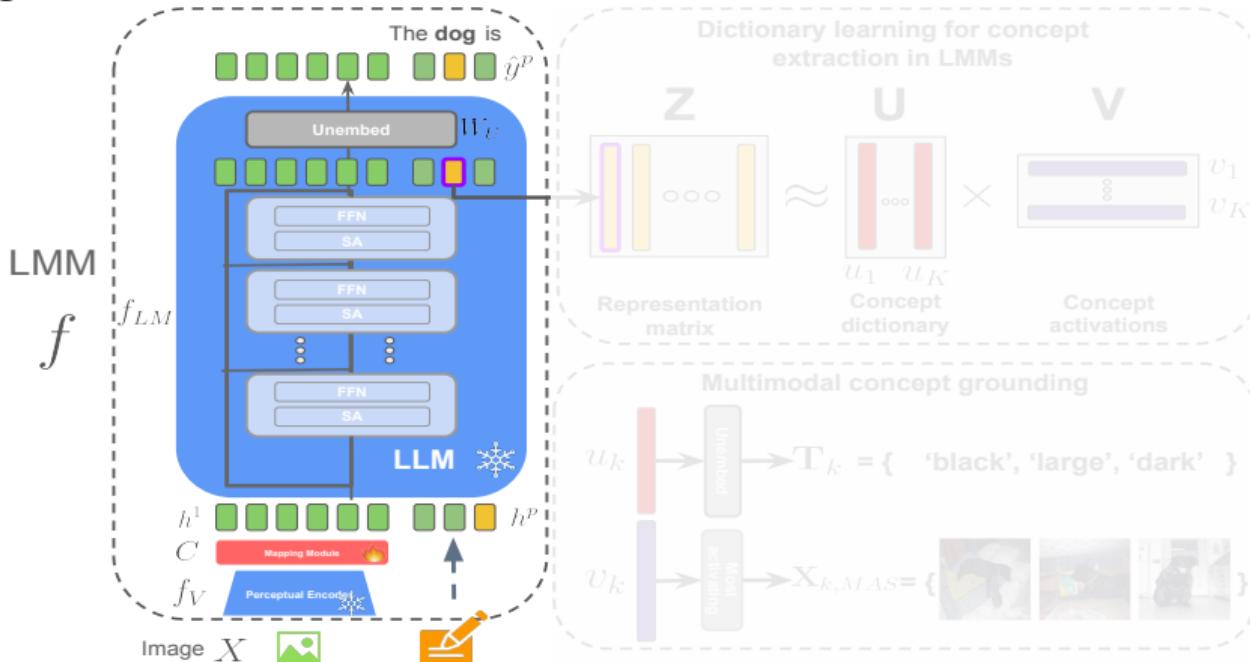


Explaining/Monitoring LMMs

For token of interest t 'Train', can we provide a multimodal concept analysis? such as:

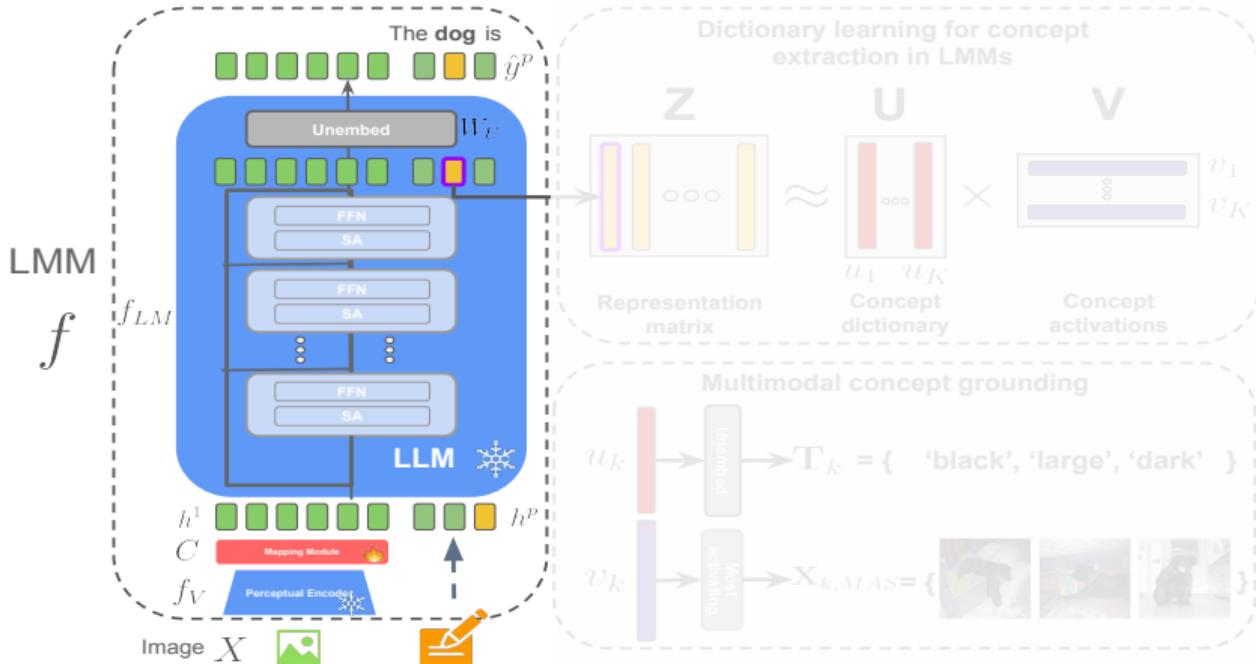


Monitoring LMM



- ▶ Input to f_{LM} - Concatenated sequence of tokens: (1) Visual tokens $C(f_V(X))$, (2) textual tokens previously predicted by f_{LM}
- ▶ Caption predicted by f_{LM} trained for next-token prediction task

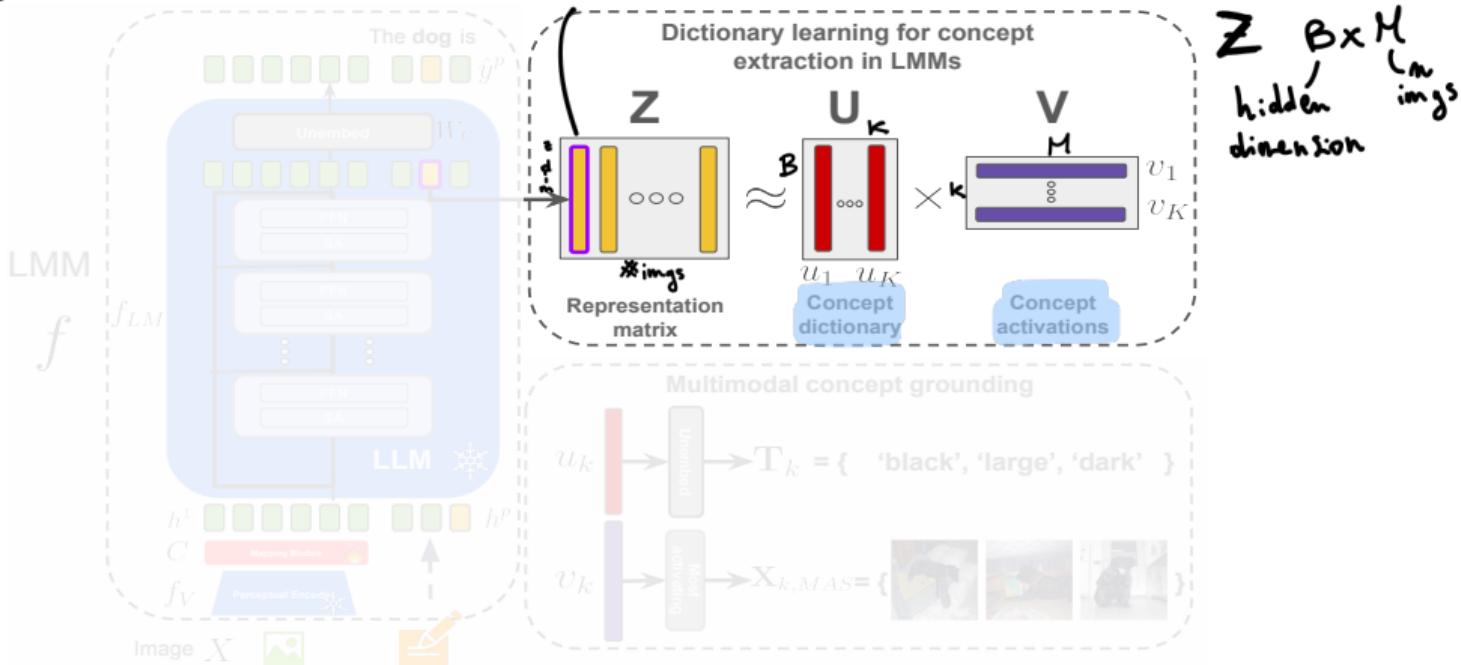
Monitoring LMM



- ▶ Extract residual stream representations of t from f for a relevant set of M images \mathbf{X}
- ▶ Collect all such B -dimensional representations as columns of matrix $\mathbf{Z} \in \mathbb{R}^{B \times M}$

Monitoring LMM

each z is representation of t for image X



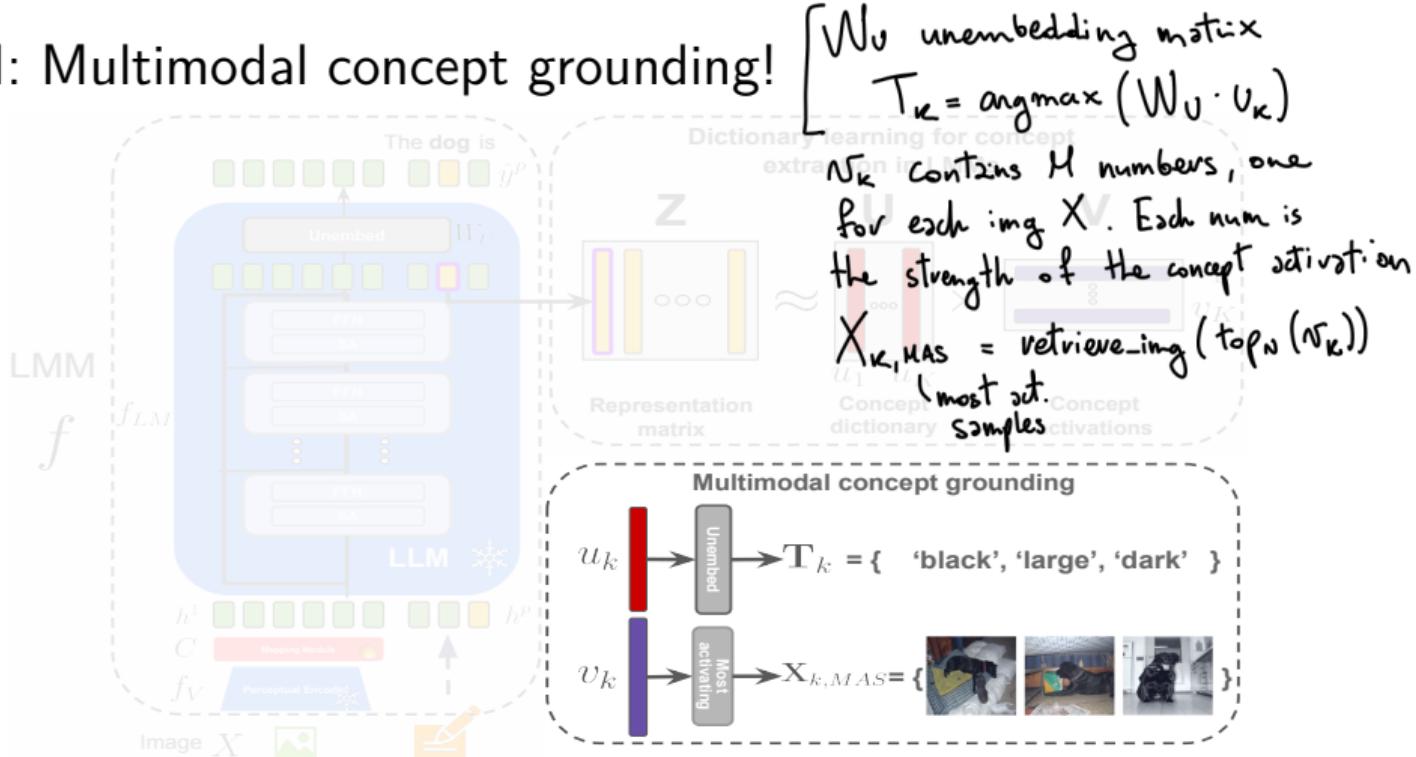
- Dictionary learning for concept extraction. Semi-NMF optimization:

$$\mathbf{U}^*, \mathbf{V}^* = \arg \min_{\mathbf{U}, \mathbf{V}} \|\mathbf{Z} - \mathbf{UV}\|_F^2 + \lambda \|\mathbf{V}\|_1 \quad s.t. \quad \mathbf{V} \geq 0, \text{ and } \|\mathbf{u}_k\|_2 \leq 1 \quad \forall k \in \{1, \dots, K\}$$

- Columns of $\mathbf{U}^* \in \mathbb{R}^{B \times K}$ – concept vectors. Rows of $\mathbf{V}^* \in \mathbb{R}^{K \times M}$ – concept activations

a L1 penalty (Lasso) → drive most v_{ik} to 0, use the fewest v_{ik}
Explanation are sparse

CoX-LMM: Multimodal concept grounding!



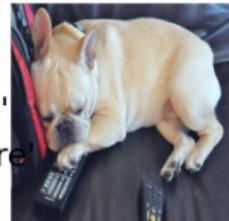
- **Text grounding:** Decode concept vector u_k with f_{LM} head and extract top tokens
- **Visual grounding:** Extract most activating samples for u_k (via activations v_k)

Example multimodal concepts

Multimodal concepts: $u_k \in \mathbf{U}^*$ simultaneously grounded in both vision and text!

- ▶ Visual: Most activating images of u_k from \mathbf{X} (via $v_k \in \mathbb{R}^M$) $\rightarrow \mathbf{X}_{k,MAS}$
- ▶ Textual: unembedding matrix W_U decode u_k and extract the most probable tokens $\rightarrow \mathbf{T}_k$

'small'



'tiny'



'puppy'



'miniature'



'cute'



'furry'



'hairy'



'fluffy'



'long'



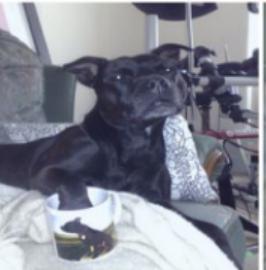
'fuzzy'



Example multimodal concepts

Multimodal concepts: $u_k \in \mathbf{U}^*$ simultaneously grounded in both vision and text!

'black'

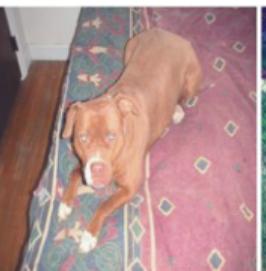
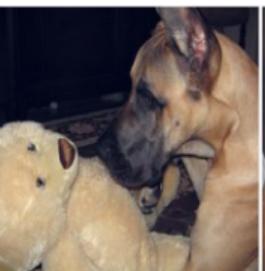
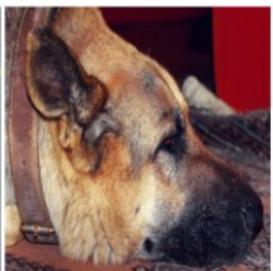


'large'

'dark'

'big'

'close'



'brown'

'large'

'dog'

'tan'

'golden'

Example multimodal concepts

Multimodal concepts: $u_k \in \mathbf{U}^*$ simultaneously grounded in both vision and text!



'dog'



'running'



'black'



'play'



'grass'



Example multimodal concepts

Multimodal concepts: $u_k \in \mathbf{U}^*$ simultaneously grounded in both vision and text!

'cat'



'kitten'



'tiger'



'rabbit'



'dog'



'herd'



'sheep'



'flock'



'farm'



'shepherd'



Example multimodal concepts

Multimodal concepts: $u_k \in \mathbf{U}^*$ simultaneously grounded in both vision and text!

'dog'



'sausage'



'hot'



'sandwich'



'plate'



Example multimodal concepts

Multimodal concepts: $u_k \in \mathbf{U}^*$ simultaneously grounded in both vision and text!



'cat'
'kitten'
'bearish'
'colored'
'pant'



'orange'
'yellow'
'cat'
'golden'
'ginger'



'kitten'
'small'
'tiny'
'baby'
'young'

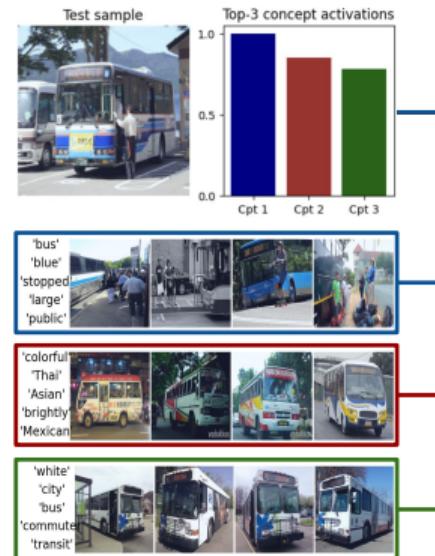
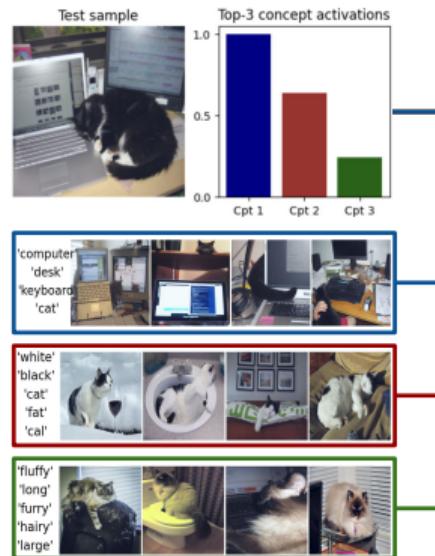
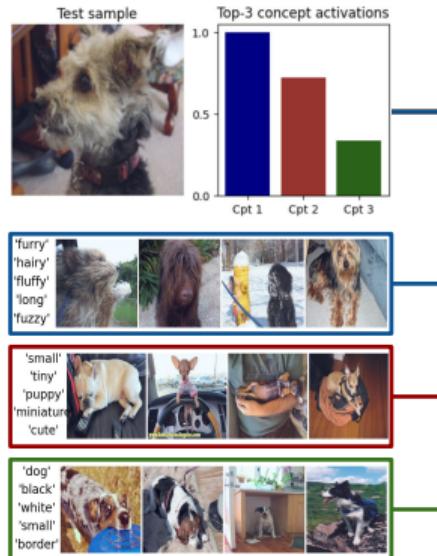


'cat'
'pet'
'black'
'dog'
'pillow'



Using the concept dictionary

- For a new image X where $t \in f(X)$, extract z_X and compute the projection on \mathbf{U}^* ,
 $v(X) = \arg \min_{v \geq 0} \|z_X - \mathbf{U}^* v\|_2^2 + \lambda \|v\|_1$
- **Most activating concepts:** From $v(X)$ we can extract the concept activations with largest magnitudes, $\tilde{u}(X)$



Using the concept dictionary

What happens if we fine-tune the LMM? It destroys previous U, V building new ones

- ▶ How do concepts encoded with the initial model change when we fine-tune it? *From generic to specific*
- ▶ Is it possible to manipulate the output of an LMM without fine-tuning it?

Yes, instead of changing W , just manipulate \tilde{v} for a certain concept u

