# Turning Polysemantic Neurons Into Pure Features by Identifying Relevant Circuits

*Maximilian Dreyer, Erblina Purelku, Johanna Vielhaben, Wojciech Samek, Sebastian Lapuschkin*

**PURE** — PURIFYING REPRESENTATIONS

Fraunhofer HHI · TECHNISCHE UNIVERSITÄT BERLIN · BIFOLD
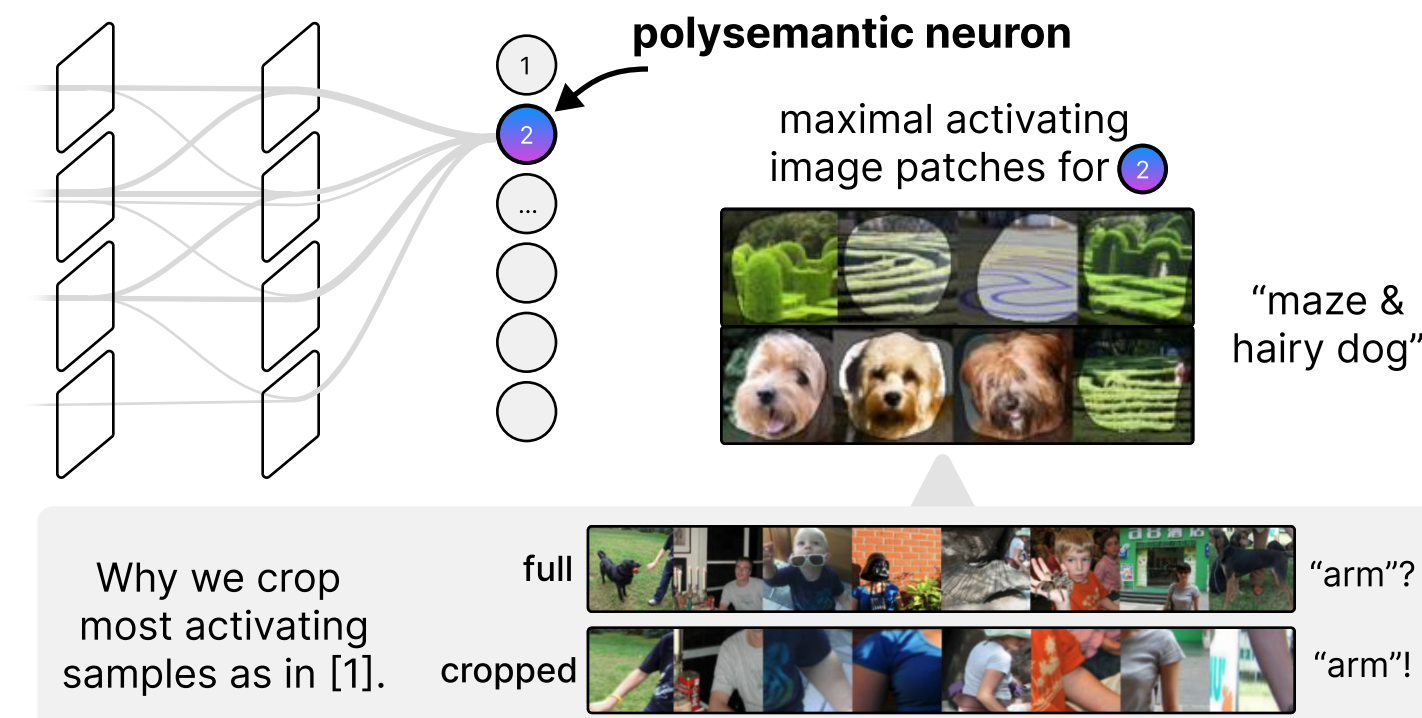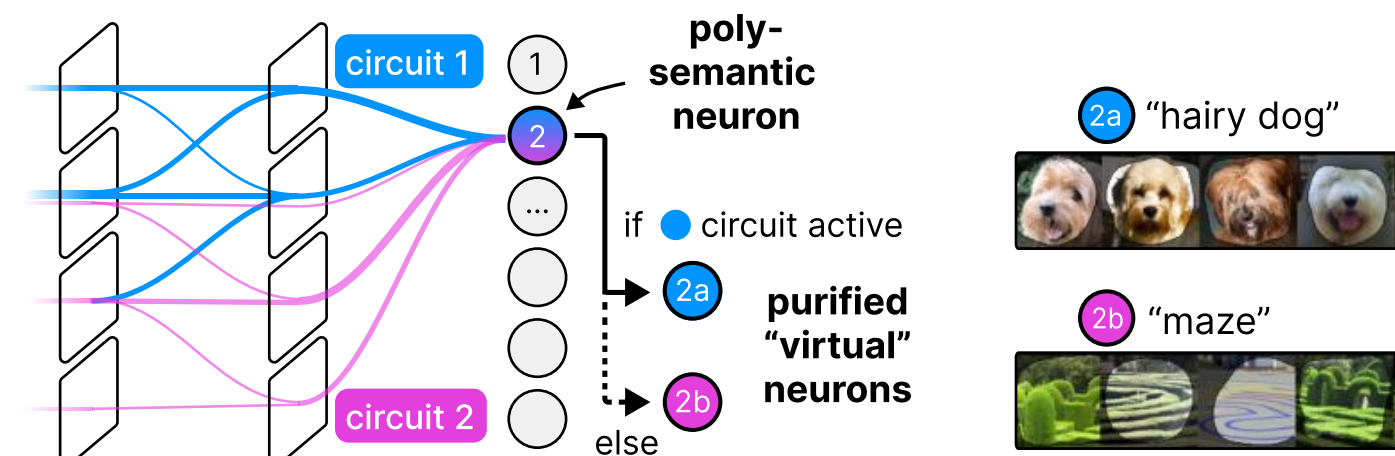
CVPR JUNE 17-21, 2024 · SEATTLE, WA

## What do neurons encode?

Studying neurons can be difficult due to polysemanticity, redundancies, etc.
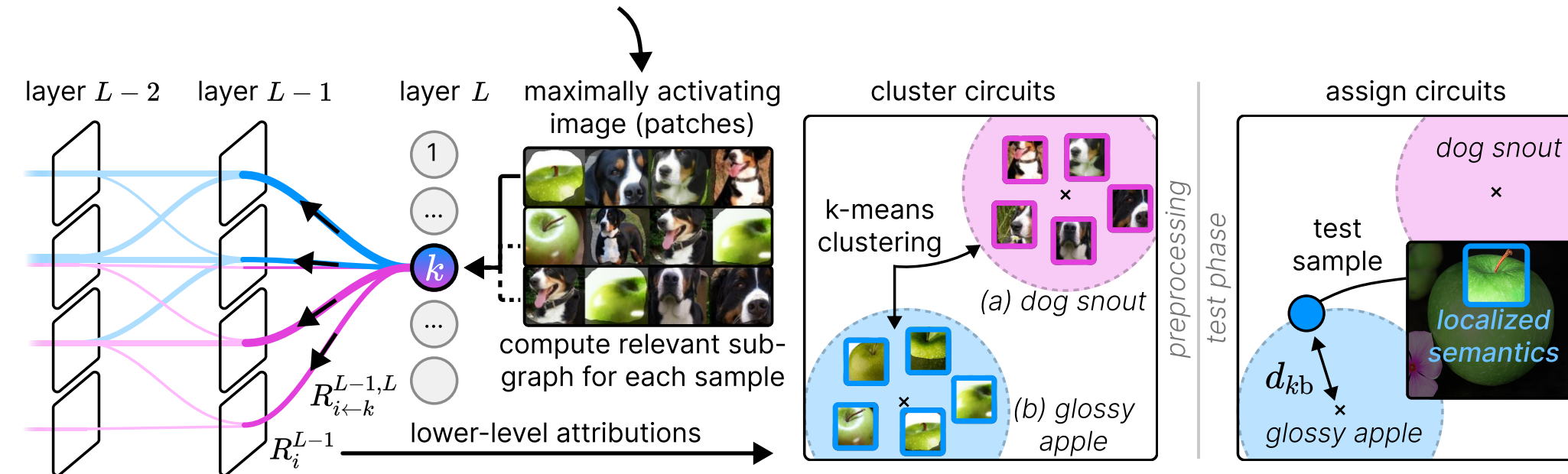
with PURE, we tackle **polysemanticity**:

**polysemantic neuron**

maximal activating image patches for ②

"maze & hairy dog"

Why we crop most activating samples as in [1].

full — "arm"?
cropped — "arm"!

## Idea

Each pure feature corresponds to a specific sub-graph.

circuit 1
poly-semantic neuron
circuit 2

if ● circuit active → **purified "virtual" neurons**
2a "hairy dog"
2b "maze"
else

When we know which sub-graph is active, we also know which feature is present.

## PURE: Purifying Representations

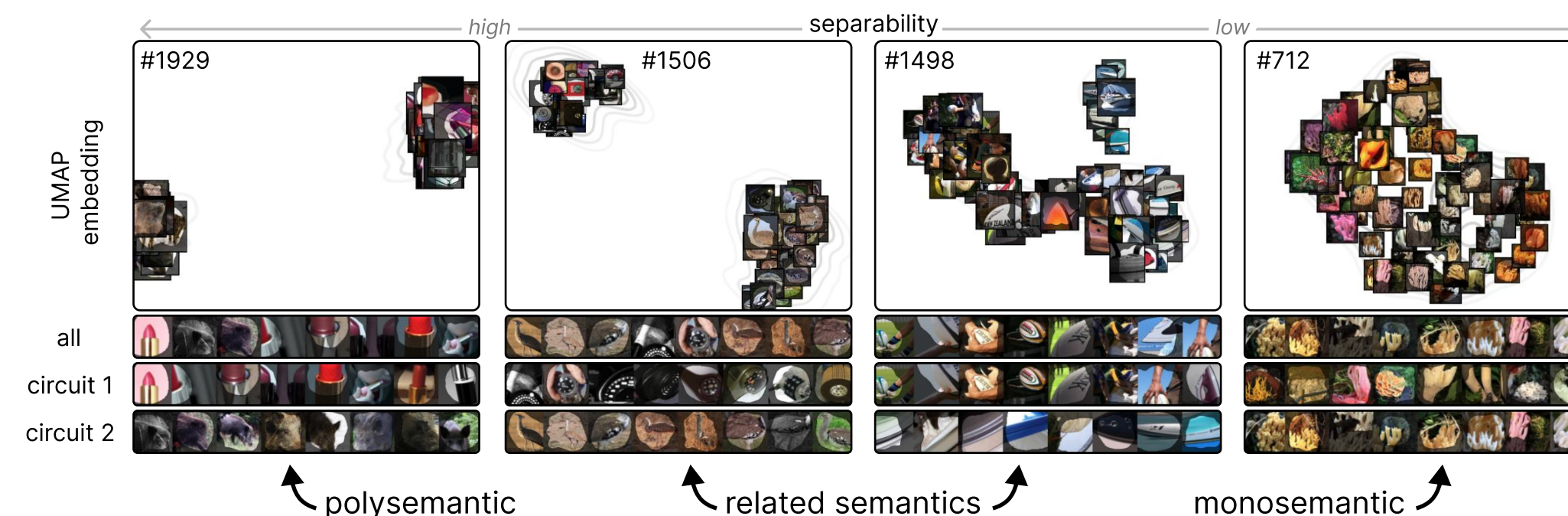1. Find **most activating samples** for a polysemantic neuron.

layer $L-2$   layer $L-1$   layer $L$

maximally activating image (patches)

compute relevant sub-graph for each sample

$R_{i \leftarrow k}^{L-1,L}$

$R_i^{L-1}$

lower-level attributions

2. **Explain neuron** and attribute lower-level neurons.

cluster circuits

k-means clustering

(a) dog snout

(b) glossy apple

*preprocessing*

assign circuits

*test phase*

dog snout

test sample

*localized semantics*

$d_{kb}$

glossy apple

3. Cluster attributions and **find circuits**.

## Qualitative Experiments

Apply PURE and sort neurons according to the effect of disentanglement.

*high* — separability — *low*

#1929    #1506    #1498    #712

UMAP embedding

all
circuit 1
circuit 2

polysemantic → related semantics ← monosemantic

## Quantitative Experiments

We use foundation model embeddings (e.g., CLIP [2] and DINOv2 [3]) to measure monosemanticity before and after purification of ResNet models.

Idea: embedding distances for maximally activating patches should decrease.

$$\mathbf{D}_{ij}^k = \sqrt{\left(\mathbf{e}_i^{\mathrm{CLIP}} - \mathbf{e}_j^{\mathrm{CLIP}}\right)^2}$$

with embedding $\mathbf{e}_i^{\mathrm{CLIP}}$ of max. act. image patch $i$ of neuron $k$

PURE achieves better diesentanglement compared to activation-based clustering.

**Monosemanticity of Clusters**

distance of CLIP embeddings

inter clusters / intra clusters

6.65   6.56   6.34   6.20 *overall*
5.87   5.94   6.09

DINOv2   PURE   activations

**Correlation of Embedding Distances**

correlation with CLIP distances

0.57   0.38   0.17

DINOv2   PURE   activations

PURE is more neuron-specific, as activations take into account all present features.

**Background Influences Activations**

relevant features for neuron #1028

*similar background*
**low activation distance** ✓

vs

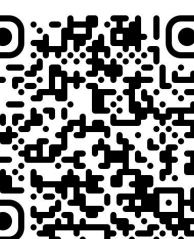*different background*
**high activation distance** ✗

## Outlook & Conclusion

→ Application to language, e.g., Large Language Models.
→ Studying the benefits of PURE for concept-based explanations, probing, and unlearning.
→ Performing an ablation study & user study for validation.

## References

[1] Achtibat, Reduan, et al. "From attribution maps to human-understandable explanations through concept relevance propagation." Nature Machine Intelligence 5.9 (2023): 1006-1019.
[2] Radford, Alec, et al. "Learning transferable visual models from natural language supervision." International conference on machine learning. PMLR, 2021.
[3] Oquab, Maxime, et al. "DINOv2: Learning Robust Visual Features without Supervision." Transactions on Machine Learning Research (2023).

PAPER   CODE