

# **CTRL-O: Language-Controllable Object-Centric Visual Representation Learning**

Aniket Didolkar\*, Andrii Zadaianchuk\*, Rabiul Awal\*, Maximilian Seitzer, Efstratios Gavves, Aishwarya Agrawal \*denotes equal contribution

### **Object-Centric Learning**

### **Pros**:

- Image decomposition into a set of slots where each slots represents a unique object.
- Works well on real-world images.

### **Cons:**

- Cannot be guided by user queries.
- Always recovers a fixed decomposition of a scene.



• CTRL-O transforms visual features into language-controllable object-centric representations. Enables precise binding between language descriptions and specific objects in complex scenes.

					Objec	t Di
Slot Init.	GT Masks	CL	DC	Binding Hits	FG-ARI	mBO
1	1	x	x	71.2	69.8	35.4
1	×	X	X	8.1	34.52	22.42
1	×	X	1	10.11	43.83	25.76
1	×	1	×	56.3	44.8	27.3
1	X	1	1	61.3	47.5	27.2

Ablation over various CTRL-O components for achieving strong visual grounding.

Approach	FG-ARI mBO		
Unsup.	DINOSAUR (MLP Dec.) [38]	40.5	27.7
	DINOSAUR (TF. Dec.) [38]	34.1	31.6
	Stable-LSD [19]	35.0	30.4
	SlotDiffusion [43]	37.3	31.4
Weak Sup.	Stable-LSD (Bbox Supervision) [39]	-	30.3
	CTRL-O (Trained on COCO)	47.5	27.2





With CTRL-O, we can control the granularity at which objects are discovered in the image.





### d Visual Grounding





The man in the grey sweatshirt



The man in the long-sleeved shirt









# **Downstream Applications: Image Generation & VQA**

images into a coherent scene.

by enabling coupling.

# Université M de Montréal

