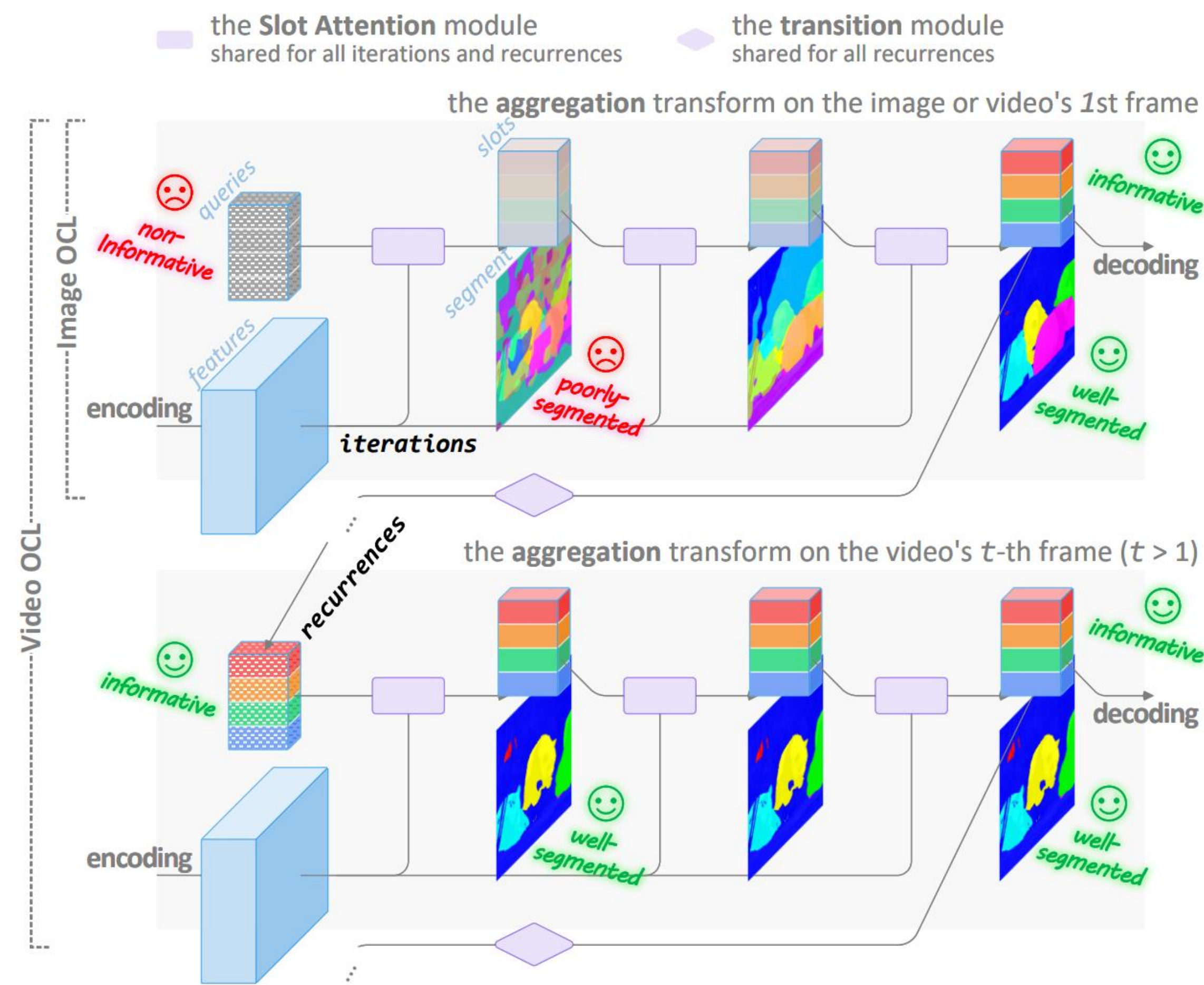


What is Object-Centric Learning (OCL):

- discover objects in images or videos without supervision;
- represent them as corresponding feature vectors, slots.

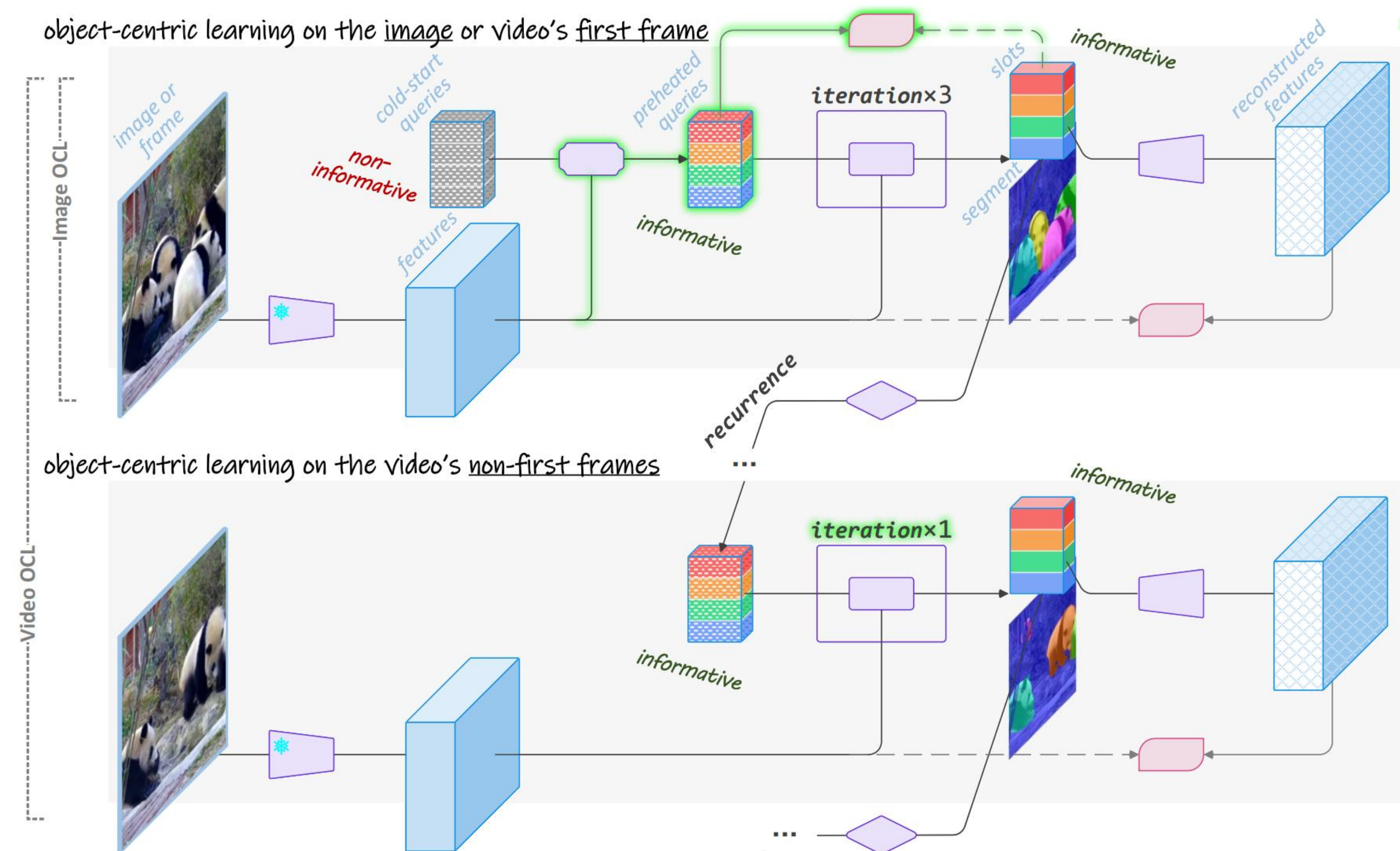
How is an OCL model trained:

- aggregate input's dense features into slots;
- reconstruct the input from such slots;
- minimize the reconstruction error.

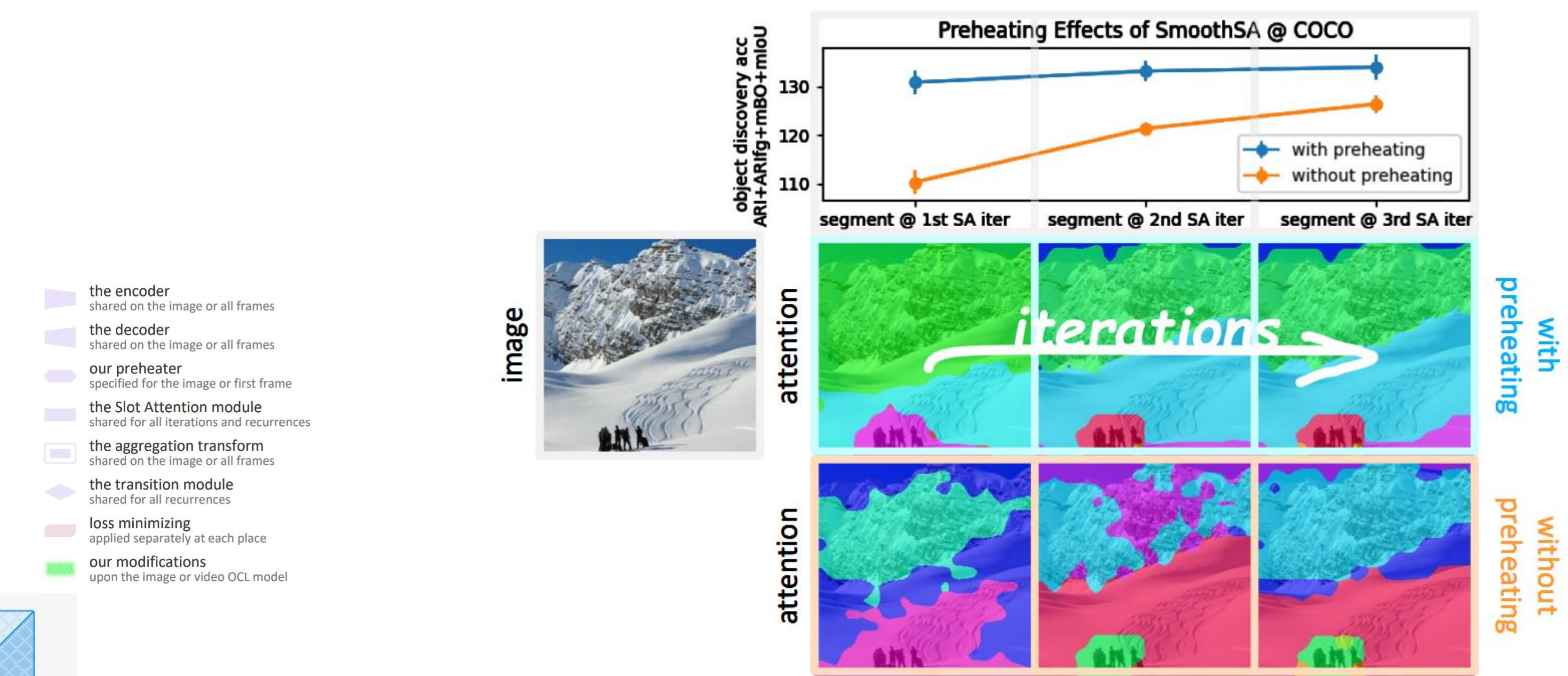


We address these issues with our SmoothSA:

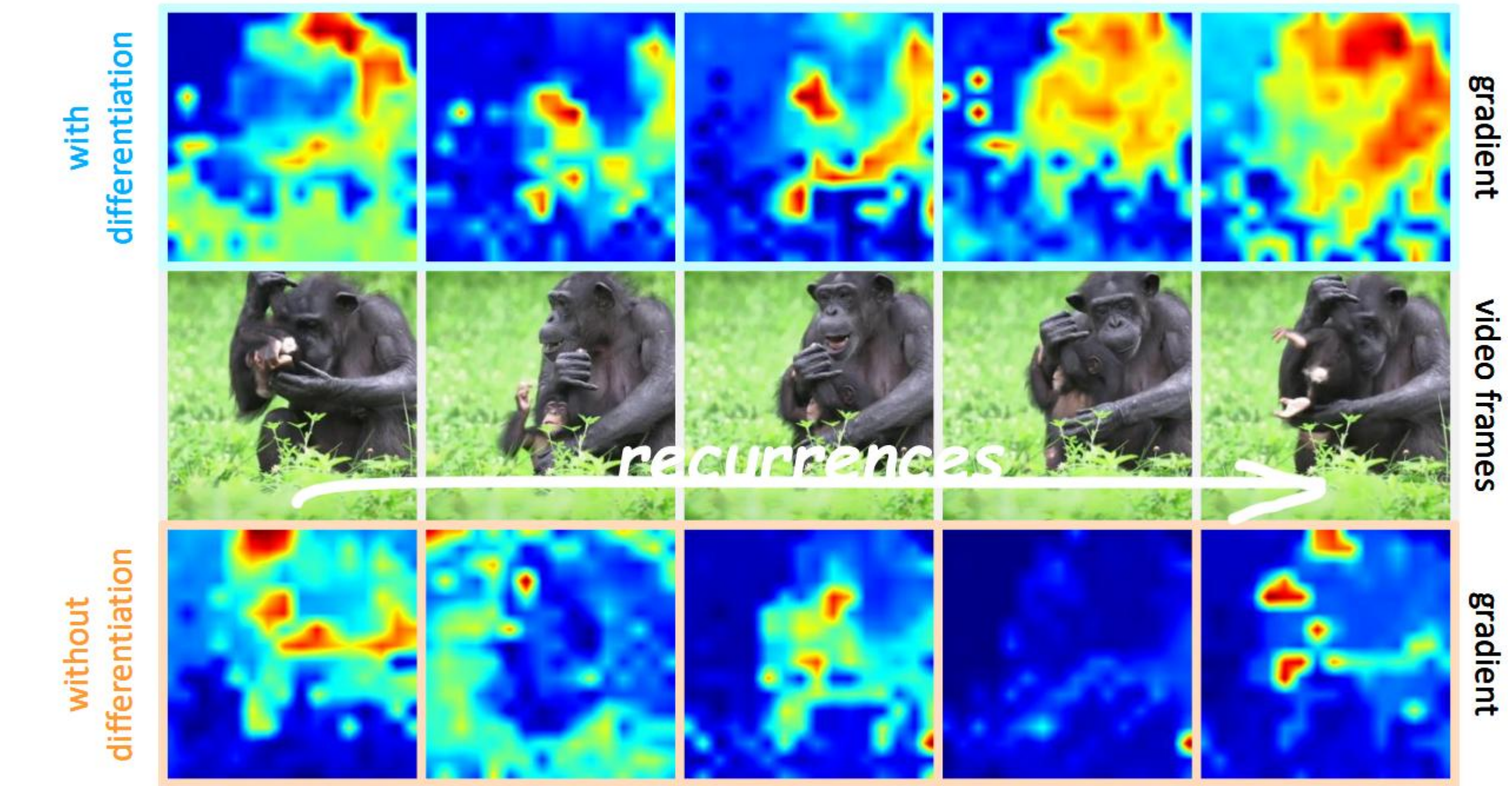
- To smooth SA iterations on image or video's first frame, we *preheat cold-start* queries with rich input-feature information, by a tiny module self-distilled inside OCL;
- To smooth SA recurrences across video's first and non-first frames, we *differentiate the homogeneous* aggregation transforms by using full and single iterations respectively



Query preheating smooths SA iterations:



Transform differentiation smooths SA recurrences:



Slot Attention (SA) in mainstream OCL:

- Image features can be aggregated into object-level representations by SA *iteratively* refining cold-start query slots.
- For video, such aggregation proceeds by SA *recurrently* shared across frames, with queries cold-started on the first frame while transitioned from the previous frame's slots thereafter.

Issues in current SA iteration and recurrence:

- *Cold-start* queries lack sample-specific cues thus hindering precise aggregation on image or video's first frame;
- Non-first frame queries are already sample-specific thus requiring aggregation transforms *different* from the first frame

SmoothSA boosts performance effectively:

	ClevrTex #slot=11				COCO #slot=7				VOC #slot=6			
	ARI	ARI _{fg}	mBO	mIoU	ARI	ARI _{fg}	mBO	mIoU	ARI	ARI _{fg}	mBO	mIoU
SLATE	17.4 _{±2.9}	87.4 _{±1.7}	44.5 _{±2.2}	43.3 _{±2.4}	17.5 _{±0.6}	28.8 _{±0.3}	26.8 _{±0.3}	25.4 _{±0.3}	18.6 _{±0.1}	26.2 _{±0.3}	37.2 _{±0.5}	36.1 _{±0.4}
DINOSAUR	50.7 _{±24.1}	89.4 _{±0.3}	53.3 _{±5.0}	52.8 _{±5.2}	18.2 _{±1.0}	37.0 _{±1.2}	28.3 _{±0.5}	26.9 _{±0.5}	21.5 _{±0.7}	36.2 _{±1.3}	40.6 _{±0.6}	39.7 _{±0.6}
SlotDiffusion	66.1 _{±1.3}	82.7 _{±1.6}	54.3 _{±0.5}	53.4 _{±0.8}	17.7 _{±0.5}	29.0 _{±0.1}	27.0 _{±0.4}	25.6 _{±0.9}	17.0 _{±1.2}	21.7 _{±1.8}	35.2 _{±0.9}	34.0 _{±1.0}
SPOT	25.6 _{±1.4}	77.1 _{±0.5}	48.2 _{±0.6}	46.3 _{±0.7}	23.7 _{±0.5}	40.4 _{±0.5}	30.9 _{±0.2}	29.3 _{±0.2}	24.5 _{±0.3}	31.0 _{±0.8}	40.1 _{±0.2}	38.6 _{±0.3}
DIAS	80.9 _{±0.3}	79.1 _{±0.3}	63.3 _{±0.1}	61.9 _{±0.0}	25.6 _{±0.1}	41.2 _{±0.3}	31.7 _{±0.1}	30.2 _{±0.1}	30.9 _{±0.5}	33.5 _{±0.7}	43.4 _{±0.5}	42.4 _{±0.5}
SmoothSA^v	76.8 _{±0.7}	82.2 _{±1.7}	60.6 _{±0.4}	58.9 _{±0.6}	29.3_{±1.0}	41.3_{±1.2}	33.4_{±0.2}	31.8_{±0.2}	35.0_{±0.3}	33.6_{±1.2}	45.2_{±0.3}	43.9_{±0.3}

Table 1. Object discovery on images. Input resolution is 224×224; DINO2 ViT-S/14 is for encoding.

	MOVi-C #slot=11, conditional				MOVi-D #slot=21, conditional				YTVIS-HQ #slot=7			
	ARI	ARI _{fg}	mBO	mIoU	ARI	ARI _{fg}	mBO	mIoU	ARI	ARI _{fg}	mBO	mIoU
STEVE	-	-	-	-	17.5 _{±0.6}	28.8 _{±0.3}	26.8 _{±0.3}	25.4 _{±0.3}	-	-	-	-
VideoSAUR	41.9 _{±1.1}	53.3 _{±2.1}	16.1 _{±0.4}	14.8 _{±0.4}	22.5 _{±5.0}	40.0 _{±0.1}	11.6 _{±0.6}	10.8 _{±1.1}	33.8 _{±0.7}	49.2 _{±0.5}	29.9 _{±0.4}	29.7 _{±0.4}
SlotContrast	64.6 _{±0.4}	59.9 _{±5.3}	27.7 _{±1.0}	25.8 _{±2.9}	45.3 _{±4.1}	63.9 _{±0.2}	26.7 _{±1.0}	25.1 _{±1.0}	37.2 _{±0.6}	49.4 _{±1.1}	33.0 _{±0.2}	32.8 _{±0.1}
RandSFQ	65.4 _{±0.7}	67.4 _{±2.1}	29.2 _{±3.8}	26.8 _{±3.7}	41.6 _{±3.7}	77.5 _{±1.0}	27.4 _{±1.0}	25.6 _{±1.0}	40.1 _{±0.4}	58.0 _{±1.0}	37.6 _{±0.4}	37.2 _{±0.4}
SmoothSA^v	50.9 _{±1.6}	69.0_{±0.3}	31.7_{±0.8}	30.2_{±0.8}	43.8 _{±1.5}	70.5 _{±0.7}	31.4_{±0.4}	30.2_{±0.4}	42.4_{±0.8}	63.0_{±3.4}	38.9_{±0.7}	38.3_{±0.6}

update 20260425 below

	MOVi-C #slot=11, conditional				MOVi-E #slot=24, conditional				YTVIS-2022 #slot=7			
	ARI	ARI _{fg}	mBO	mIoU	ARI	ARI _{fg}	mBO	mIoU	ARI	ARI _{fg}	mBO	mIoU
VideoSAUR	41.9 _{±1.1}	53.3 _{±2.1}	16.1 _{±0.4}	14.8 _{±0.4}	17.4 _{±2.5}	34.6 _{±0.7}	8.3 _{±4.9}	7.5 _{±4.3}	33.4 _{±0.8}	48.2 _{±0.7}	27.2 _{±0.3}	26.8 _{±0.3}
SlotContrast	64.6 _{±0.4}	59.9 _{±5.3}	27.7 _{±1.0}	25.8 _{±2.9}	29.9 _{±4.9}	70.6 _{±3.8}	20.7 _{±1.4}	19.3 _{±1.2}	35.2 _{±0.8}	51.4 _{±0.7}	29.7 _{±0.5}	29.3 _{±0.6}
RandSFQ	65.4 _{±0.7}	67.4 _{±2.1}	29.2 _{±3.8}	26.8 _{±3.7}	30.5 _{±1.2}	82.1 _{±3.1}	23.0 _{±1.2}	21.6 _{±1.4}	37.9 _{±1.3}	51.8 _{±1.2}	32.2 _{±1.8}	31.5 _{±1.8}
SmoothSA^v	50.9 _{±1.6}	69.0_{±0.3}	31.7_{±0.8}	30.2_{±0.8}	36.7_{±0.6}	73.6 _{±0.6}	28.6_{±0.1}	27.4_{±0.1}	42.0_{±0.6}	59.0_{±2.1}	36.0_{±0.5}	34.9_{±0.6}

Table 2. Object discovery on videos. Input resolution is 224×224; DINO2 ViT-S/14 is for encoding.

		COCO #slot=7			
		class top1	top3	bbox IoU	#match
SPOT	+ MLP	88.7 _{±0.3}	97.2 _{±0.1}	48.6 _{±0.1}	5061 _{±30}
SmoothSA^v	+ MLP	89.3_{±0.5}	97.3_{±0.2}	49.5_{±1.1}	5210_{±223}

update 20260425 below

		YTVIS-2022 #slot=7			
		class top1	top3	bbox IoU	#match
SlotContrast+MLP		87.1 _{±0.2}	96.4 _{±0.1}	48.2 _{±0.3}	19943 _{±156}
SmoothSA^v	+ MLP	91.0_{±0.2}	97.6_{±0.1}	42.6_{±1.4}	8957_{±34}

Table 3. Object recognition on images and videos.

		GQA #slot=7	
		accuracy %	
SPOT	+ Aloe	52.3 _{±2.8}	
SmoothSA^v	+ Aloe	56.7_{±1.9}	

update 20260425 below

		CLEVRER #slot=7	
		per option %	per question %
SlotContrast	+ Aloe	97.2 _{±1.1}	95.6 _{±0.9}
SmoothSA^v	+ Aloe	98.7_{±0.4}	96.9_{±0.6}

Table 4. Visual question answering on images and videos.

