

# 000 SMOOTHING SLOT ATTENTION ITERATIONS AND RE- 001 CURRENCES 002 003 004

005 **Anonymous authors**

006 Paper under double-blind review

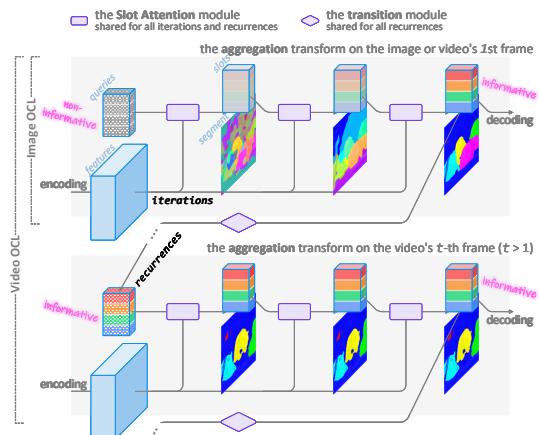
## 007 ABSTRACT 008 009

011 Slot Attention (SA) and its variants lie at the heart of mainstream Object-Centric  
012 Learning (OCL). Objects in an image can be aggregated into corresponding slot  
013 vectors, by *iteratively* refining cold-start query vectors, typically three times, via  
014 SA on image features. For video, this aggregation is *recurrently* shared across  
015 frames, with queries cold-started on the first frame while transitioned from the pre-  
016 vious frame’s slots on non-first frames. However, cold-start queries lack sample-  
017 specific cues thus hindering precise aggregation on the image or video’s first  
018 frame; Also, non-first frames’ queries are already sample-specific thus requiring  
019 aggregation transforms different from the first frame. We address these issues  
020 for the first time with our *SmoothSA*: (1) To smooth SA iterations on the image  
021 or video’s first frame, we *preheat* the cold-start queries with rich information of  
022 input features, via a tiny module self-distilled inside OCL; (2) To smooth SA re-  
023 currances across all video frames, we *differentiate* the homogeneous transforms on  
024 the first and non-first frames, by using full and single iterations respectively. Com-  
025 prehensive experiments on object discovery, recognition and downstream bench-  
026 marks validate our method’s effectiveness. Further analyses illuminate how our  
027 method smooths SA iterations and recurrences. Our source code and training logs  
028 are provided in the supplement.

## 029 1 INTRODUCTION 030

031 Object-Centric Learning (OCL) (Locatello  
032 et al., 2020) aims to represent objects in a vi-  
033 sual scene as distinct vectors, with the back-  
034 ground as another vector. Ideally, this yields a  
035 structured compact representation that outper-  
036 forms popular dense feature maps in advanced  
037 vision tasks. In dynamics modeling, evolving  
038 these object-level slots over time captures more  
039 accurate object interactions (Villar-Corrales &  
040 Behnke, 2025). For visual reasoning, their con-  
041 cise form allows more explicit object rela-  
042 tionship modeling, slashing the search space and  
043 computation load (Ding et al., 2021). In vi-  
044 sual prediction, disentangling objects facilitates  
045 more compositional generation of future frames  
(Villar-Corrales et al., 2023).

046 Powered by Slot Attention (SA) (Locatello  
047 et al., 2020), modern OCL methods have sig-  
048 nificantly improved and can now scale to real-  
049 world complex images and videos. SA is essen-  
050 tially a form of iterative cross attention, where  
051 query vectors compete to aggregate their cor-  
052 responding object information, discovering ob-  
053 jects as segmentation masks and representing  
them as slot vectors (Locatello et al., 2020).



054 Figure 1: Image Object-Centric Learning (OCL) is re-  
055 alized via Slot Attention (SA) *iterations* on the image,  
056 while video OCL is via SA *recurrences* across video  
057 frames. In SA iterations on the image or video’s first  
058 frame, the **cold-start queries** lack information for ac-  
059 curate aggregation; In SA recurrences across video’s first  
060 and non-first frames, the **homogeneous trans-  
061 forms**, i.e., the fixed three SA iterations, cannot jointly  
062 adapt to the first and non-first queries, which have a sig-  
063 nificant information gap.

054 The model is trained by minimizing reconstruction loss based on the slots, requiring no external  
 055 supervision. Specifically, for image, the queries are usually cold-start and sampled from multiple  
 056 Gaussian distributions fitted to the entire dataset (Jia et al., 2023). Such queries contain no informa-  
 057 tion about any specific sample, thus to obtain slots by refining queries using SA on image features,  
 058 typically three iterations are necessary. For video, such aggregation occurs recurrently across all  
 059 frames in a shared way, where queries for the first frame are the same as in the image case while  
 060 queries for non-first frames are transitioned from the previous frame’s slots (Singh et al., 2022b).  
 061 Unlike the first frame’s queries, non-first frames’ queries are already quite sample-specific. But the  
 062 aggregation transforms are identical or homogeneous across all frames.

063 To the best of our knowledge, all works on SA and its variants confront these facts but have not  
 064 acknowledged the implied issues, as shown in Figure 1: (i1) *Query cold-start* in SA iterations. For  
 065 an image or video’s first frame, the cold-start queries lack scene-specific information. Although  
 066 three SA iterations can gradually refine these non-informative queries into useful slots, such aggrega-  
 067 tion would not work as good as that with informative queries. (i2) *Transform homogeneity* in SA  
 068 recurrences. For video frames, the first frame’s queries are cold-start while non-first frames’ are  
 069 much more informative. These differing conditions impose different requirements on the aggrega-  
 070 tion transforms, thus such homogeneous transforms would not work as good as those adapted to  
 071 informative-different queries.

072 Our solution is simple yet effective. We propose *SmoothSA*, which smooths SA iterations on the  
 073 image or video’s first frame by preheating the queries, and smooths SA recurrences across video’s  
 074 first and non-first frames by differentiating the transforms: (s1) A tiny module *preheats* the cold-start  
 075 queries using rich information from input features. It is trained by predicting current slots through  
 076 self-distillation within the OCL model. (s2) Different aggregation transforms handle video’s first  
 077 and non-first frames respectively. This is realized by simply employing three SA iterations on the  
 078 first frame while only one on each non-first frame.

079 Briefly, our contributions are: (c1) for the first time addressing the query cold-start issue in SA  
 080 iterations on the image and video’s first frame; (c2) for the first time addressing the transform ho-  
 081 mogeneity issue in SA recurrences across the first and non-first frames; (c3) new state-of-the-art  
 082 on both image and video OCL benchmarks; (c4) consistent performance boosts on downstream  
 083 advanced vision tasks.

## 085 2 RELATED WORK

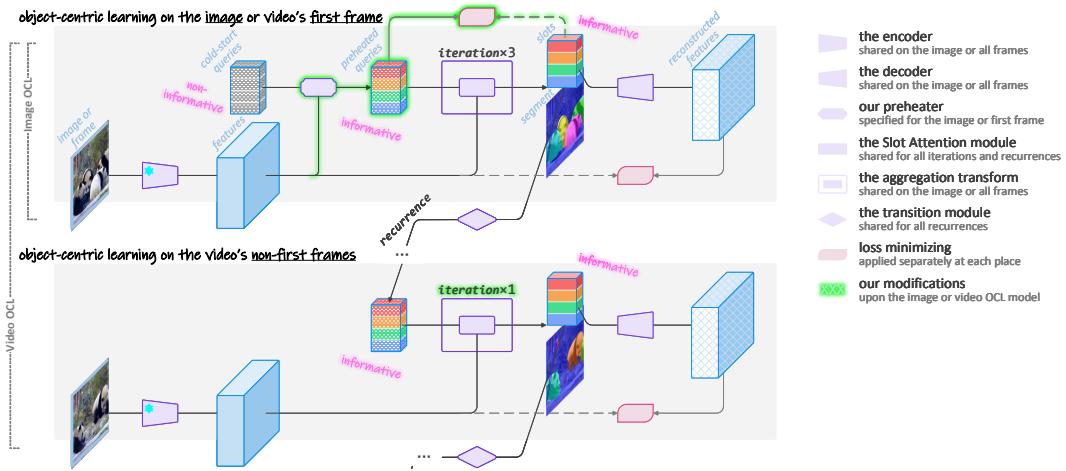
087 As SA is a kind of cross attention that depends on queries to aggregate information from visual  
 088 features, we review works from perspectives of aggregation and queries.

089 **Slot Attention on images and videos.** The seminal work on the aggregation module SA (Locatello  
 090 et al., 2020) proposes refining the initial randomly initialized queries into object-centric slots via  
 091 typically three iterations of the same SA module on image features. Then, all image OCL methods  
 092 including (Singh et al., 2022a; Seitzer et al., 2023; Wu et al., 2023b; Jiang et al., 2023; Kakogeorgiou  
 093 et al., 2024; Zhao et al., 2025b;c;d;e) adopt this iterative design. The pioneering work STEVE  
 094 (Singh et al., 2022b) extends SA to videos by conducting standard image OCL on each frame,  
 095 using randomly initialized queries for the first frame while using recurrently predicted queries from  
 096 previous slots for non-first frames. After, all video OCL methods including SAVi (Kipf et al., 2022),  
 097 SAVi++ (Elsayed et al., 2022), SOLV (Aydemir et al., 2023), VideoSAUR (Zadaianchuk et al.,  
 098 2024), SlotContrast (Manasyan et al., 2025), STATM (Li et al., 2025b), SlotPi (Li et al., 2025a) and  
 099 RandSF.Q (Zhao et al., 2025a) adopt such recurrent design. Now that SA is the core module of  
 100 mainstream OCL methods for images or videos, all methods face but never acknowledge two issues  
 101 described above. Our method is the first to address these issues directly.

102 **Query initialization for Slot Attention iterations.** For images, the initial queries serve as the  
 103 starting point for aggregation based on SA iterations. The principal contradiction is that no object  
 104 cues are available before aggregation. SA (Locatello et al., 2020) initializes queries by drawing  
 105 multiple samples from a global Gaussian distribution, which is learned on the entire dataset and  
 106 embeds global cues for object discovery. BO-QSA (Jia et al., 2023) proposes learning multiple  
 107 Gaussian distributions so that more distinct cues are embedded into initial queries, thus enabling  
 108 better aggregation. However, the queries are still cold-start. MetaSlot (Liu et al., 2025) takes two

108 steps: firstly initializing queries from multiple Gaussians for draft aggregation iterations, and then  
 109 replacing the draft slots with object embeddings from a large codebook (Van Den Oord et al., 2017)  
 110 for additional aggregation iterations. This mitigates the iterative query cold-start effectively, but still  
 111 relies on cold-start queries. We directly address such iterative query cold-start issue.

112 **Query prediction for Slot Attention recurrences.** For video’s first frame, the queries can be ob-  
 113 tained in the same way as in the image case, or by transforming cues like object bounding boxes  
 114 in SAVi (Kipf et al., 2022) and SAVi++ (Elsayed et al., 2022), albeit at the cost of extra expensive  
 115 annotations. For non-first frames, the queries are predicted from the previous frame’s slots. STEVE  
 116 (Singh et al., 2022b) and most other OCL methods use a Transformer encoder block for such re-  
 117 current prediction. STADM (Li et al., 2025b) and SlotPi (Li et al., 2025a) employ auto-regressive  
 118 Transformer encoder variants for the same purpose. The most recent work RandSFQ (Zhao et al.,  
 119 2025a) incorporates the next frame’s feature for more informative query prediction, and uses random  
 120 slot-feature pairs for explicit query prediction learning, significantly boosting OCL performance on  
 121 videos. However, improving query prediction alone will never reach the core issue, recurrent trans-  
 122 form discrepancy. We directly address this recurrent transform homogeneity issue.



139 Figure 2: The overall model and our modifications. *(upper)* For image OCL, we preheat the cold-start queries  
 140 to be informative so as to smooth SA iterations on the image (or video’s first frame). Our preheater is a tiny  
 141 module trained to predict vectors approximating the slots as the preheated queries from the cold-start queries  
 142 and image features. *(upper + lower)* For video OCL, we differentiate the homogeneous transforms to adapt to  
 143 the first and non-first queries, non-informative and informative respectively, to smooth SA recurrences across  
 144 all frames. This is achieved by using three SA iterations on the first frame and one on non-first frames.

### 3 PROPOSED METHOD

149 Mainstream image or video OCL methods confront two issues: the query cold-start in SA iterations  
 150 on the image or video’s first frame, and the transform homogeneity in SA recurrences across video’s  
 151 first and non-first frames. We address these issues for the first time with our *SmoothSA*, by preheating  
 152 queries to smooth SA iterations and differentiating transforms to smooth SA recurrences.

#### 3.1 SLOT ATTENTION ITERATION AND RECURRENCE

153 Mainstream OCL methods mainly take the encoder-aggregator-decoder model design (Zhao et al.,  
 154 2025d): The encoder encodes the image or video frames into features, the aggregator aggregates  
 155 features into slots, and the decoder decodes slots into the reconstruction of the input in some form as  
 156 the source of supervision. The aggregator, which is based on Slot Attention (SA) (Locatello et al.,  
 157 2020) or its variants, is the core of OCL, so let us focus on it.

158 **SA iterations on the image or video’s first frame.** An SA-based aggregator  $\phi_a$  takes cold-start  
 159 vectors  $Q_1 \in \mathbb{R}^{n \times c}$  as the query, and input features  $F_1 \in \mathbb{R}^{h \times w \times c}$  as the key and value.  $\phi_a$  is

162 applied on the query, key and value typically three times, to refine the query iteratively into object-  
 163 level feature vectors  $\mathbf{S}_1 \in \mathbb{R}^{n \times c}$ , i.e., slots, as the sparse representation of the visual scene:

$$164 \quad \mathbf{Q}_1 = \phi_n(\mathbf{C}) \quad (1)$$

$$167 \quad \mathbf{S}_1, \mathbf{M}_1 = \Phi_a(\mathbf{Q}_1, \mathbf{F}_1) \quad (2)$$

168 where the aggregation transform  $\Phi_a$  can be expanded into:

$$170 \quad \mathbf{S}_1^{(0)} := \mathbf{Q}_1 \quad (2a)$$

$$171 \quad \mathbf{S}_1^{(i)}, \mathbf{M}_1^{(i)} = \phi_a(\mathbf{S}_1^{(i-1)}, \mathbf{F}_1) \quad i = 1, 2, 3 \quad (2b)$$

$$173 \quad \mathbf{S}_1, \mathbf{M}_1 := \mathbf{S}_1^{(3)}, \mathbf{M}_1^{(3)} \quad (2c)$$

174 In Equation (1), if cues  $\mathbf{C}$  are  $n$  slots to use, then the initializer  $\phi_n$  samples  $n$  vectors as the queries  
 175  $\mathbf{Q}_1$  from its trainable Gaussian distribution(s) (Locatello et al., 2020; Jia et al., 2023); If cues  $\mathbf{C}$   
 176 are the bounding boxes of objects in the video’s first frame, then the initializer  $\phi_n$  projects cues  $\mathbf{C}$   
 177 into the queries  $\mathbf{Q}_1$  (Kipf et al., 2022; Elsayed et al., 2022). In whichever case, queries  $\mathbf{Q}_1$  lack  
 178 sample-specific information, namely, being cold-start.

179 Considering that  $\mathbf{F}_1$  is the high-quality feature of the image or video’s first frame, typically extracted  
 180 by vision foundation model DINO2 (Oquab et al., 2023), the quality of the transform  $\Phi_a$  is decided  
 181 by queries  $\mathbf{Q}_1$ . Therefore, if we could preheat the cold-start queries  $\mathbf{Q}_1$  to be more informative, the  
 182 aggregation transform  $\Phi_a$  on the image or video’s first frame would perform better.

183 **SA recurrences across video’s first and non-first frames.** The transform  $\Phi_a$  based on SA iterations  
 184 is shared across all frames recurrently. Namely, the transform  $\Phi_a$  happens across both first and  
 185 non-first frames, where the former is identical to the image case formulated in Equations (1) and (2).  
 186 For the latter, queries  $\mathbf{Q}_t$  are recurrently transitioned from previous frame’s slots  $\mathbf{S}_{t-1}$ :

$$188 \quad \mathbf{Q}_t = \phi_r(\mathbf{S}_{t-1}) \quad t \geq 2 \quad (3)$$

$$191 \quad \mathbf{S}_t, \mathbf{M}_t = \Phi'_a(\mathbf{Q}_t, \mathbf{F}_t) \quad (4)$$

192 where the aggregation transform  $\Phi'_a$  can be expanded into:

$$193 \quad \mathbf{S}_t^{(0)} := \mathbf{Q}_t \quad (4a)$$

$$195 \quad \mathbf{S}_t^{(i)}, \mathbf{M}_t^{(i)} = \phi_a(\mathbf{S}_t^{(i-1)}, \mathbf{F}_t) \quad i = 1, 2, 3 \quad (4b)$$

$$196 \quad \mathbf{S}_t, \mathbf{M}_t := \mathbf{S}_t^{(3)}, \mathbf{M}_t^{(3)} \quad (4c)$$

198 In Equation (3), the transitioner  $\phi_r$  takes previous frame’s slots  $\mathbf{S}_{t-1}$  as input and predicts current  
 199 queries  $\mathbf{Q}_t$ . Considering that  $\mathbf{S}_{t-1}$  is the information-intensive representation of the previous frame  
 200 and that the transitioner  $\phi_r$  learns knowledge of transition dynamics (Singh et al., 2022b), current  
 201 queries  $\mathbf{Q}_t$  is actually informative to current frame. This is different from the first frame queries  $\mathbf{Q}_1$ ,  
 202 which is cold-start and thus non-informative.

203 The non-first frames’ transform shares exactly the same SA module from the first transform with  
 204 the same number of SA iterations, i.e.,  $\Phi'_a \equiv \Phi_a$ . On the other hand, the information gap between  
 205 the first queries  $\mathbf{Q}_1$  and non-first queries  $\mathbf{Q}_t$  imposes different requirements on these transforms.  
 206 Therefore, if we could differentiate the homogeneous transforms  $\Phi_a$  and  $\Phi'_a$  for the first and non-  
 207 first frames respectively, the aggregation across the first and non-first frames would be better.

### 208 3.2 PREHEATING COLD-START QUERIES

211 To overcome the query cold-start issue and smooth SA iterations on the image or video’s first frame,  
 212 we preheat the cold-start queries with rich information from input features. A tiny module is trained  
 213 via self-distillation inside the OCL model, to predict vectors that approximate the aggregated slots  
 214 as the preheated queries, from the cold-start queries conditioned on input features.

215 Our chain-of-thought is as follows: (i) Informative slots can be aggregated by iteratively refining  
 216 non-informative queries; (ii) More informative queries contribute to better slots aggregation; (iii)

216 How to preheat the queries to be more informative? (iv) Aligning the preheated queries with the  
 217 aggregated slots, which are quite informative.

218 Firstly, we insert this between Equations (1) and (2):

$$220 \quad \mathbf{Q}_1^* = \phi_p(\mathbf{Q}_1, \mathbf{F}_1) \quad (5)$$

221 where the preheater  $\phi_p$  is parameterized as a single Transformer decoder block (Vaswani et al.,  
 222 2017), whose self-attention and cross-attention are switched. This is because exchanging information  
 223 among non-informative queries firstly is meaningless. Please refer to Table 5 ablation studies  
 224 for why not using an extra SA module as the preheater, and for why switching the self-attention and  
 225 cross-attention.

226 Secondly, we replace Equation (2a) with:

$$228 \quad \mathbf{S}_1^{(0)} := \text{sg}(\mathbf{Q}_1^*) \quad (6)$$

229 where  $\text{sg}(\cdot)$  is stopping gradient. Stopping gradient flow from the SA module  $\phi_a$  to the preheated  
 230 queries  $\mathbf{Q}_1^*$  disentangles the training of  $\phi_a$  and  $\phi_p$ . Please refer to Table 5 ablation studies for why  
 231 stopping gradient flow on the preheated queries.

232 Lastly, to obtain the preheating ability, we train our preheater  $\phi_p$  with the following objective:

$$234 \quad \arg \min_{C, \phi_n, \phi_p} \text{MSE}(\mathbf{Q}_1^*, \text{sg}(\mathbf{S}_1)) \quad (7)$$

235 where the MSE loss is combined with the original OCL loss(es). To ensure the sufficient training of  
 236  $\phi_p$ , we can use a relatively large coefficient for it. Please refer to Table 5 ablation studies for what  
 237 weight to set for such preheating loss.

238 **Comment 1.** Our preheater is trained with OCL intermediate results as the ground-truth, without  
 239 any external supervision, forming rigid self-distillation. This is also bootstrap, as good slots  $\mathbf{S}_1$   
 240 leads to better preheated queries  $\mathbf{Q}_1^*$ , and in turn better  $\mathbf{Q}_1^*$  leads to better  $\mathbf{S}_1$ .

241 **Comment 2.** Our preheater is similar to the SA module, without an heavy RNN module. Thus our  
 242 preheater introduces approximately less than 1/3 more computation overhead in the aggregation,  
 243 which is negligible considering the heavy computation of encoding and decoding.

### 244 3.3 DIFFERENTIATING HOMOGENEOUS TRANSFORMS

245 To overcome the transform homogeneity issue and smooth SA recurrences across the video's first  
 246 and non-first frames, we differentiate the homogeneous transforms for the first and non-first frames  
 247 respectively. For the different transform requirements due to the gap between the first cold-start  
 248 queries and non-first informative queries, full and single SA iterations are used respectively.

249 Our chain-of-thought is: (i) First frame queries are non-informative, thus three SA iterations are  
 250 needed to refine the queries into good slots; (ii) Non-first frame queries are already informative, thus  
 251 a single SA iteration is enough.

252 As mentioned above, the first-frame transform  $\Phi_a$  and non-first frame transforms  $\Phi'_a$  are identical  
 253 in all existing methods but should be different. There are two ways to differentiate them: (1) use  
 254 separate SA parameters for  $\Phi_a$  and  $\Phi'_a$ ; (2) use different number of iterations for  $\Phi_a$  and  $\Phi'_a$ .  
 255 We choose the second solution. This is because  $\Phi_a$  and  $\Phi'_a$  should learn the general aggregation  
 256 capability in each SA iteration and sharing enforces this. Please refer to Table 5 ablation studies for  
 257 what numbers of iterations for first and non-first transforms to set.

258 We simply reduce the number of SA iterations in non-first frame transforms  $\Phi'_a$  to once, while  
 259 always use three SA iterations in the first frame transform  $\Phi_a$ . Namely, we keep Equations (2b)  
 260 and (2c) unchanged, while replacing Equations (4b) and (4c) with:

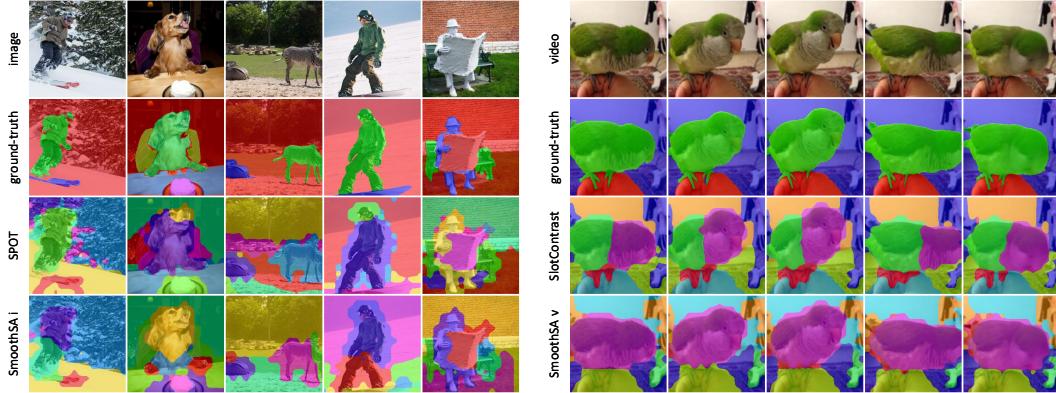
$$261 \quad \mathbf{S}_t^{(i)}, \mathbf{M}_t^{(i)} = \phi_a(\mathbf{S}_t^{(i-1)}, \mathbf{F}_t) \quad i = 1 \quad (8b)$$

$$262 \quad \mathbf{S}_t, \mathbf{M}_t := \mathbf{S}_t^{(1)}, \mathbf{M}_t^{(1)} \quad (8c)$$

263 For conditional SA like in SAVi (Kipf et al., 2022) and SAVi++ (Elsayed et al., 2022), they use  
 264 homogeneous aggregation transforms, consisting of one single SA iteration for all frames. But we

270 still use three SA iterations on the first frame and one on non-first frames. They believe that objects' 271 bounding boxes as query initialization is informative enough. But in fact, they still carry little object 272 information, except the spatial information. Thus more iterations on the first frame is still necessary. 273 Their ablation study leads them to believe that one iteration is better than more just because they 274 were not aware of such recurrent transform homogeneity issue. Please refer to Table 5 ablation 275 studies for what numbers of iterations for first and non-first transforms to set.

276 **Comment 3.** Our differentiated transforms have 2/3 less computation overhead in aggregation on 277 non-first frames, which is negligible considering the heavy encoding and decoding.



293 Figure 3: Qualitative results of our SmoothSA on images (*left*) and videos (*right*), compared with SotA  
294 methods SPOT and SlotContrast respectively.

## 297 4 EXPERIMENT

299 We conduct experiments on object discovery along with downstream tasks, object recognition and  
300 visual question answering, to evaluate our slots representation quality, with three random seeds.  
301

### 303 4.1 INSTANTIATING SMOOTHSA

305 As shown in Figure 2, our OCL model with SmoothSA is based on DIAS<sup>i</sup> (Zhao et al., 2025e) for  
306 images and RandSF.Q (Zhao et al., 2025a) for videos, respectively. These two state-of-the-art (SotA)  
307 methods share identical designs except techniques specific to image and video. For image OCL, we  
308 remove slots pruning tricks from DIAS<sup>i</sup>, and then replace its SA variant with our SmoothSA. For  
309 video OCL, we use RandSF.Q as it is, and then replace its SA with our SmoothSA. Thus we have  
310 models SmoothSA<sup>i</sup> and SmoothSA<sup>v</sup>, where *i* is image and *v* is video.

311 Note that for conditional video OCL like SAVi (Kipf et al., 2022) and SAVi++ (Elsayed et al., 2022),  
312 the authors always use one SA iteration on all frames. But whether it is conditional or not, we always  
313 use three SA iterations on the first frame while one iteration on non-first frames.

	ClevrTex #slot=11				COCO #slot=7				VOC #slot=6			
	ARI	ARI <sub>fg</sub>	mBO	mIoU	ARI	ARI <sub>fg</sub>	mBO	mIoU	ARI	ARI <sub>fg</sub>	mBO	mIoU
SLATE	17.4 <sub>±2.9</sub>	87.4 <sub>±1.7</sub>	44.5 <sub>±2.2</sub>	43.3 <sub>±2.4</sub>	17.5 <sub>±0.6</sub>	28.8 <sub>±0.3</sub>	26.8 <sub>±0.3</sub>	25.4 <sub>±0.3</sub>	18.6 <sub>±0.1</sub>	26.2 <sub>±0.8</sub>	37.2 <sub>±0.5</sub>	36.1 <sub>±0.4</sub>
DINOSAUR	50.7 <sub>±24.1</sub>	89.4 <sub>±0.3</sub>	53.3 <sub>±5.0</sub>	52.8 <sub>±5.2</sub>	18.2 <sub>±1.0</sub>	37.0 <sub>±1.2</sub>	28.3 <sub>±0.5</sub>	26.9 <sub>±0.5</sub>	21.5 <sub>±0.7</sub>	36.2 <sub>±1.3</sub>	40.6 <sub>±0.6</sub>	39.7 <sub>±0.6</sub>
SlotDiffusion	66.1 <sub>±1.3</sub>	82.7 <sub>±1.6</sub>	54.3 <sub>±0.5</sub>	53.4 <sub>±0.8</sub>	17.7 <sub>±0.5</sub>	29.0 <sub>±0.1</sub>	27.0 <sub>±0.4</sub>	25.6 <sub>±0.4</sub>	17.0 <sub>±1.2</sub>	21.7 <sub>±1.8</sub>	35.2 <sub>±0.9</sub>	34.0 <sub>±1.0</sub>
SPOT	25.6 <sub>±1.4</sub>	77.1 <sub>±0.7</sub>	48.3 <sub>±0.5</sub>	46.4 <sub>±0.6</sub>	20.0 <sub>±0.5</sub>	40.0 <sub>±0.7</sub>	30.2 <sub>±0.3</sub>	28.6 <sub>±0.3</sub>	20.3 <sub>±0.7</sub>	33.5 <sub>±1.1</sub>	40.1 <sub>±0.5</sub>	38.7 <sub>±0.7</sub>
DIAS <sup>i</sup>	80.9 <sub>±0.3</sub>	79.1 <sub>±0.3</sub>	63.3 <sub>±0.0</sub>	61.9 <sub>±0.0</sub>	22.0 <sub>±0.2</sub>	41.4 <sub>±0.2</sub>	31.1 <sub>±0.1</sub>	29.7 <sub>±0.1</sub>	26.6 <sub>±1.0</sub>	33.7 <sub>±1.5</sub>	43.3 <sub>±0.3</sub>	42.4 <sub>±0.3</sub>
SmoothSA <sup>i</sup>	76.8 <sub>±1.4</sub>	80.8 <sub>±1.6</sub>	60.0 <sub>±1.8</sub>	58.1 <sub>±2.2</sub>	26.2 <sub>±0.8</sub>	42.1 <sub>±0.7</sub>	33.2 <sub>±0.4</sub>	31.7 <sub>±0.4</sub>	30.6 <sub>±0.6</sub>	34.3 <sub>±0.5</sub>	45.3 <sub>±0.5</sub>	44.1 <sub>±0.6</sub>

322 Table 1: Object discovery on images. Input resolution is 224×224; DINO2 ViT-S/14 is for encoding.  
323

324 4.2 OBJECT DISCOVERY  
325

326 In mainstream OCL methods, attention maps of the slots are binarized as the byproduct object  
327 segmentation, i.e., discovering objects. This intuitively reflects slots' representation quality.  
328

329 On image datasets ClevrTex<sup>1</sup>, COCO<sup>2</sup> and VOC<sup>3</sup>, we  
330 compare our SmoothSA<sup>i</sup> with baselines SLATE (Singh  
331 et al., 2022a), DINOSAUR (Seitzer et al., 2023), SlotD-  
332 iffusion (Wu et al., 2023b), SPOT (Kakogeorgiou et al.,  
333 2024) (no distillation and finetuning tricks) and DIAS  
334 (Zhao et al., 2025e) (no slot pruning). On video dataset  
335 YouTube Video Instance Segmentation<sup>4</sup> (YTVIS) the  
336 high-quality version<sup>5</sup>, we compare our SmoothSA<sup>v</sup> with  
337 baselines STEVE (Singh et al., 2022b), VideoSAUR  
338 (Zadaianchuk et al., 2024), SlotContrast (Manasyan  
339 et al., 2025) and RandSF.Q (Zhao et al., 2025a). The  
340 performance metrics are ARI<sup>6</sup>, ARI<sub>fg</sub> (foreground),  
341 mBO (Uijlings et al., 2013) and mIoU<sup>7</sup>. ARI score is  
342 calculated with the segmentation area as the weight, thus ARI mainly reflects how well the back-  
343 ground is segmented while ARI<sub>fg</sub> reflects how well large objects are segmented. mBO shows how  
344 objects that are best overlapped with the ground-truth are segmented. mIoU is the most strick metric.  
345 **Note that, unless otherwise specified, we use image ARI, ARI<sub>fg</sub>, mBO and mIoU for object discov-  
346 ery on images, while using video ones for object discovery on videos. Also note that on dataset  
YTVIS, we use video clip length 5 for training while 20 for evaluation.**  
347

348 As shown in Table 1, on synthetic dataset ClevrTex, our SmoothSA<sup>i</sup> is as competitive as the latest  
349 SotA DIAS<sup>i</sup> and significantly better than former SotA SPOT in all metrics. On real-world dataset  
350 COCO, our SmoothSA<sup>i</sup> is consistently better than DIAS<sup>i</sup> in all metrics, 4+ points in ARI. On real-  
351 world dataset VOC, our method pushes the ARI value forward by 4 points. Our method achieves  
352 overall new SotA in ARI, mBO and mIoU, except relative limited performance boosts in ARI<sub>fg</sub>.  
353

354 As shown in Table 2, on real-world video dataset YTVIS, our SmoothSA<sup>v</sup> defeats all baselines by a  
355 large margin, even including the latest super SotA method RandSF.Q, which has already pushed the  
356 older SotA performance significantly forward by up to 10 points.  
357

358 4.3 OBJECT RECOGNITION  
359

360 Besides the byproduct segmentation, recognizing the discovered objects' attributes like class and  
361 bounding box from the slots can directly reflect the object-centric representation quality.  
362

363 On real-world image dataset COCO, we compare  
364 our SmoothSA<sup>i</sup> with baseline SPOT (Kakogeorgiou  
365 et al., 2024). On real-world video dataset YTVIS,  
366 we compare our SmoothSA<sup>v</sup> with baseline SlotCon-  
367 trast (Manasyan et al., 2025). We follow the routine  
368 of (Seitzer et al., 2023): firstly convert all images  
369 into slots representation, with some threshold filter-  
370 ing; then train a two-layer MLP model to classify and  
371 regress the matched object's class label and bound-  
372 ing box coordinates in a supervised way. We use top1  
373

	ARI	ARI <sub>fg</sub>	mBO	mIoU
YTVIS #slot=7, #step=20				
VideoSAUR	34.6 <sub>±0.5</sub>	48.6 <sub>±0.7</sub>	31.4 <sub>±0.3</sub>	31.2 <sub>±0.3</sub>
SlotContrast	38.7 <sub>±0.9</sub>	48.9 <sub>±0.7</sub>	35.0 <sub>±0.3</sub>	34.9 <sub>±0.3</sub>
DIAS <sup>v</sup>	33.6 <sub>±0.4</sub>	49.3 <sub>±0.7</sub>	36.1 <sub>±1.4</sub>	35.2 <sub>±0.8</sub>
RandSF.Q	42.0 <sub>±0.3</sub>	59.4 <sub>±1.1</sub>	39.8 <sub>±0.3</sub>	39.4 <sub>±0.3</sub>
SmoothSA <sup>v</sup>	44.1 <sub>±1.8</sub>	61.5 <sub>±3.2</sub>	41.1 <sub>±1.4</sub>	40.6 <sub>±1.4</sub>

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

	class	top1↑	bbox	R2↑
COCO #slot=7				
SPOT	+	MLP	0.67 <sub>±0.0</sub>	0.62 <sub>±0.1</sub>
SmoothSA <sup>i</sup>	+	MLP	0.73 <sub>±0.0</sub>	0.64 <sub>±0.1</sub>
YTVIS #slot=7, #step=20				
SlotContrast	+	MLP	0.40 <sub>±0.1</sub>	0.53 <sub>±0.1</sub>
SmoothSA <sup>v</sup>	+	MLP	0.50 <sub>±0.0</sub>	0.62 <sub>±0.0</sub>

Table 3: Object recognition on image dataset COCO and video dataset YTVIS.

<sup>1</sup><https://www.robots.ox.ac.uk/~vgg/data/clevrte>

<sup>2</sup><https://cocodataset.org>

<sup>3</sup><http://host.robots.ox.ac.uk/pascal/VOC>

<sup>4</sup><https://youtube-vos.org/dataset/vis>

<sup>5</sup><https://github.com/SysCV/vmt?tab=readme-ov-file#hq-ytvvis-high-quality-video-instance-segmentation-dataset>

<sup>6</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.adjusted\\_rand\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.adjusted_rand_score.html)

<sup>7</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.jaccard\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.jaccard_score.html)

accuracy<sup>8</sup> to measure the classification performance, and R2 score<sup>9</sup> to measure the regression performance.

As shown in Table 3, the object recognition accuracy on both real-world complex images and videos are improved a lot by using our method as the slots representation extractor, compared with that using baseline methods. This demonstrates the high quality of our slots representation.

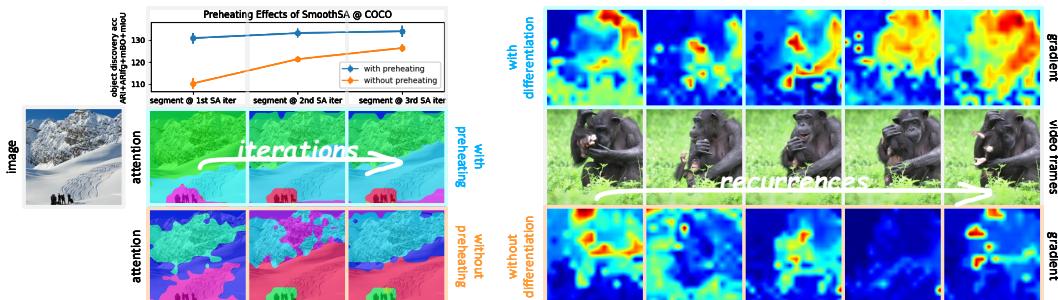
#### 384 4.4 VISUAL QUESTION ANSWERING

386 In visual question answering (VQA) tasks, the visual modality slots are combined with language  
387 modality words embeddings together, testing the representation quality and versatility further.

389 For VQA on images, we compare our  
390 SmoothSA<sup>i</sup> plus multi-modal reasoning model  
391 Aloe (Ding et al., 2021) with baseline SPOT  
392 plus Aloe on real-world complex image dataset  
393 GQA<sup>10</sup>. For VQA on videos, we compare  
394 our SmoothSA<sup>v</sup> plus Aloe with baseline Slot-  
395 Contrast plus Aloe on synthetic video dataset  
396 CLEVRER<sup>11</sup>. Please note that for the image  
397 dataset, we use Aloe as it is while on the  
398 video dataset we introduce temporal embed-  
399 ding scheme from (Wu et al., 2023a). For  
400 the upstream OCL models, we firstly pre-  
401 train them on corresponding datasets and freeze  
402 them to represent samples as slots. These visual input along with textual inputs representing ques-  
403 tions are fed into the Aloe model together, appended with a classification token. The output is  
404 obtained by projecting the transformed classification token into logits of all possible class labels,  
i.e., answers.

405 As shown in table 4, using our method as the upstream model improves the image VQA performance  
406 on dataset GQA by 4+ points. As for video VQA on CLEVRER, using our method as the upstream  
407 boosts the performance too, whether measured by per option accuracy or per question accuracy.

## 408 5 DISCUSSION



422 Figure 4: (left, middle row) Using query preheating, good segmentation can be obtained at the beginning  
423 of SA iterations, with even better segmentation at the end. (right, top row) Using transform differentia-  
424 tion, balanced gradient signals can be obtained across SA recurrences, even showing good object contours.

### 425 Query preheating smooths SA iterations

427 The segmentation accuracies generally increase along with the SA iterations, so we expect that our  
428 query preheating provides better initial queries and the accuracies increase faster. Please refer to Sec-

<sup>8</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html)

<sup>9</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.r2\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.r2_score.html)

<sup>10</sup><https://cs.stanford.edu/people/dorarad/gqa>

<sup>11</sup><http://clevrer.csail.mit.edu>

432  
 433  
 434  
 435  
 436  
 437  
 438  
 439  
 440  
 441  
 442  
 443  
 444  
 445  
 446  
 447  
 448  
 449  
 450  
 451  
 452  
 453  
 454  
 455  
 456  
 457  
 458  
 459  
 460  
 461  
 462  
 463  
 464  
 465  
 466  
 467  
 468  
 469  
 470  
 471  
 472  
 473  
 474  
 475  
 476  
 477  
 478  
 479  
 480  
 481  
 482  
 483  
 484  
 485  
 486  
 487  
 488  
 489  
 490  
 491  
 492  
 493  
 494  
 495  
 496  
 497  
 498  
 499  
 500  
 501  
 502  
 503  
 504  
 505  
 506  
 507  
 508  
 509  
 510  
 511  
 512  
 513  
 514  
 515  
 516  
 517  
 518  
 519  
 520  
 521  
 522  
 523  
 524  
 525  
 526  
 527  
 528  
 529  
 530  
 531  
 532  
 533  
 534  
 535  
 536  
 537  
 538  
 539  
 540  
 541  
 542  
 543  
 544  
 545  
 546  
 547  
 548  
 549  
 550  
 551  
 552  
 553  
 554  
 555  
 556  
 557  
 558  
 559  
 560  
 561  
 562  
 563  
 564  
 565  
 566  
 567  
 568  
 569  
 570  
 571  
 572  
 573  
 574  
 575  
 576  
 577  
 578  
 579  
 580  
 581  
 582  
 583  
 584  
 585  
 586  
 587  
 588  
 589  
 590  
 591  
 592  
 593  
 594  
 595  
 596  
 597  
 598  
 599  
 600  
 601  
 602  
 603  
 604  
 605  
 606  
 607  
 608  
 609  
 610  
 611  
 612  
 613  
 614  
 615  
 616  
 617  
 618  
 619  
 620  
 621  
 622  
 623  
 624  
 625  
 626  
 627  
 628  
 629  
 630  
 631  
 632  
 633  
 634  
 635  
 636  
 637  
 638  
 639  
 640  
 641  
 642  
 643  
 644  
 645  
 646  
 647  
 648  
 649  
 650  
 651  
 652  
 653  
 654  
 655  
 656  
 657  
 658  
 659  
 660  
 661  
 662  
 663  
 664  
 665  
 666  
 667  
 668  
 669  
 670  
 671  
 672  
 673  
 674  
 675  
 676  
 677  
 678  
 679  
 680  
 681  
 682  
 683  
 684  
 685  
 686  
 687  
 688  
 689  
 690  
 691  
 692  
 693  
 694  
 695  
 696  
 697  
 698  
 699  
 700  
 701  
 702  
 703  
 704  
 705  
 706  
 707  
 708  
 709  
 710  
 711  
 712  
 713  
 714  
 715  
 716  
 717  
 718  
 719  
 720  
 721  
 722  
 723  
 724  
 725  
 726  
 727  
 728  
 729  
 730  
 731  
 732  
 733  
 734  
 735  
 736  
 737  
 738  
 739  
 740  
 741  
 742  
 743  
 744  
 745  
 746  
 747  
 748  
 749  
 750  
 751  
 752  
 753  
 754  
 755  
 756  
 757  
 758  
 759  
 760  
 761  
 762  
 763  
 764  
 765  
 766  
 767  
 768  
 769  
 770  
 771  
 772  
 773  
 774  
 775  
 776  
 777  
 778  
 779  
 780  
 781  
 782  
 783  
 784  
 785  
 786  
 787  
 788  
 789  
 790  
 791  
 792  
 793  
 794  
 795  
 796  
 797  
 798  
 799  
 800  
 801  
 802  
 803  
 804  
 805  
 806  
 807  
 808  
 809  
 810  
 811  
 812  
 813  
 814  
 815  
 816  
 817  
 818  
 819  
 820  
 821  
 822  
 823  
 824  
 825  
 826  
 827  
 828  
 829  
 830  
 831  
 832  
 833  
 834  
 835  
 836  
 837  
 838  
 839  
 840  
 841  
 842  
 843  
 844  
 845  
 846  
 847  
 848  
 849  
 850  
 851  
 852  
 853  
 854  
 855  
 856  
 857  
 858  
 859  
 860  
 861  
 862  
 863  
 864  
 865  
 866  
 867  
 868  
 869  
 870  
 871  
 872  
 873  
 874  
 875  
 876  
 877  
 878  
 879  
 880  
 881  
 882  
 883  
 884  
 885  
 886  
 887  
 888  
 889  
 890  
 891  
 892  
 893  
 894  
 895  
 896  
 897  
 898  
 899  
 900  
 901  
 902  
 903  
 904  
 905  
 906  
 907  
 908  
 909  
 910  
 911  
 912  
 913  
 914  
 915  
 916  
 917  
 918  
 919  
 920  
 921  
 922  
 923  
 924  
 925  
 926  
 927  
 928  
 929  
 930  
 931  
 932  
 933  
 934  
 935  
 936  
 937  
 938  
 939  
 940  
 941  
 942  
 943  
 944  
 945  
 946  
 947  
 948  
 949  
 950  
 951  
 952  
 953  
 954  
 955  
 956  
 957  
 958  
 959  
 960  
 961  
 962  
 963  
 964  
 965  
 966  
 967  
 968  
 969  
 970  
 971  
 972  
 973  
 974  
 975  
 976  
 977  
 978  
 979  
 980  
 981  
 982  
 983  
 984  
 985  
 986  
 987  
 988  
 989  
 990  
 991  
 992  
 993  
 994  
 995  
 996  
 997  
 998  
 999  
 1000  
 1001  
 1002  
 1003  
 1004  
 1005  
 1006  
 1007  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015  
 1016  
 1017  
 1018  
 1019  
 1020  
 1021  
 1022  
 1023  
 1024  
 1025  
 1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079  
 1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133  
 1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187  
 1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241  
 1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295  
 1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349  
 1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403  
 1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457  
 1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511  
 1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565  
 1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619  
 1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673  
 1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727  
 1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781  
 1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835  
 1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889  
 1890  
 1891  
 1892  
 1893  
 1894  
 1895  
 1896  
 1897  
 1898  
 1899  
 1900  
 1901  
 1902  
 1903  
 1904  
 1905  
 1906  
 1907  
 1908  
 1909  
 1910  
 1911  
 1912  
 1913  
 1914  
 1915  
 1916  
 1917  
 1

- 486 Jindong Jiang, Fei Deng, Gautam Singh, and Sungjin Ahn. Object-Centric Slot Diffusion. *Advances*  
 487 *in Neural Information Processing Systems*, 2023.
- 488
- 489 Ioannis Kakogeorgiou, Spyros Gidaris, Konstantinos Karantzalos, and Nikos Komodakis. Spot:  
 490 Self-Training with Patch-Order Permutation for Object-Centric Learning with Autoregressive  
 491 Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
 492 *Recognition*, pp. 22776–22786, 2024.
- 493 Thomas Kipf, Gamaleldin Elsayed, Aravindh Mahendran, et al. Conditional Object-Centric Learn-  
 494 ing from Video. *International Conference on Learning Representations*, 2022.
- 495
- 496 Jian Li, Wan Han, Ning Lin, Yu-Liang Zhan, Ruizhi Chengze, Haining Wang, Yi Zhang, Hongsheng  
 497 Liu, Zidong Wang, Fan Yu, et al. SlotPi: Physics-informed Object-centric Reasoning Models.  
 498 *arXiv preprint arXiv:2506.10778*, 2025a.
- 499 Jian Li, Pu Ren, Yang Liu, and Hao Sun. Reasoning-Enhanced Object-Centric Learning for Videos.  
 500 In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining*  
 501 V. 1, pp. 659–670, 2025b.
- 502
- 503 Hongjia Liu, Rongzhen Zhao, Haohan Chen, and Joni Pajarinen. Metaslot: Break through the fixed  
 504 number of slots in object-centric learning. *arXiv preprint arXiv:2505.20772*, 2025.
- 505
- 506 Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, et al. Object-Centric Learning with  
 507 Slot Attention. *Advances in Neural Information Processing Systems*, 33:11525–11538, 2020.
- 508
- 509 Anna Manasyan, Maximilian Seitzer, Filip Radovic, Georg Martius, and Andrii Zadaianchuk. Tem-  
 510 porally consistent object-centric learning by contrasting slots. In *Proceedings of the Computer*  
*Vision and Pattern Recognition Conference*, pp. 5401–5411, 2025.
- 511
- 512 Maxime Oquab, Timothee Darcet, Theo Moutakanni, et al. DINOv2: Learning Robust Visual Fea-  
 513 tures without Supervision. *Transactions on Machine Learning Research*, 2023.
- 514
- 515 Jiyang Qi, Yan Gao, Yao Hu, Xinggang Wang, Xiaoyu Liu, Xiang Bai, Serge Belongie, Alan Yuille,  
 516 Philip HS Torr, and Song Bai. Occluded video instance segmentation: A benchmark. *Inter-  
 517 national Journal of Computer Vision*, 130(8):2022–2039, 2022.
- 518
- 519 Maximilian Seitzer, Max Horn, Andrii Zadaianchuk, et al. Bridging the Gap to Real-World Object-  
 520 Centric Learning. *International Conference on Learning Representations*, 2023.
- 521
- 522 Gautam Singh, Fei Deng, and Sungjin Ahn. Illiterate DALL-E Learns to Compose. *International*  
*523 Conference on Learning Representations*, 2022a.
- 524
- 525 Gautam Singh, Yi-Fu Wu, and Sungjin Ahn. Simple Unsupervised Object-Centric Learning for  
 526 Complex and Naturalistic Videos. *Advances in Neural Information Processing Systems*, 35:  
 527 18181–18196, 2022b.
- 528
- 529 Jasper RR Uijlings, Koen EA Van De Sande, Theo Gevers, and Arnold WM Smeulders. Selective  
 530 Search for Object Recognition. *International Journal of Computer Vision*, 104:154–171, 2013.
- 531
- 532 Aaron Van Den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural Discrete Representation  
 533 Learning. *Advances in Neural Information Processing Systems*, 30, 2017.
- 534
- 535 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
 536 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-  
 537 tion processing systems*, 30, 2017.
- 538
- 539 Angel Villar-Corrales and Sven Behnke. Playslot: Learning inverse latent dynamics for control-  
 540 lable object-centric video prediction and planning. In *Forty-second International Conference on*  
*541 Machine Learning*, 2025.
- 542
- 543 Angel Villar-Corrales, Ismail Wahdan, and Sven Behnke. Object-centric video prediction via de-  
 544 coupling of object dynamics and interactions. In *2023 IEEE International Conference on Image*  
*545 Processing (ICIP)*, pp. 570–574. IEEE, 2023.

- 540 Ziyi Wu, Nikita Dvornik, Klaus Greff, Thomas Kipf, and Animesh Garg. SlotFormer: Unsupervised  
 541 Visual Dynamics Simulation with Object-Centric Models. *International Conference on Learning*  
 542 *Representations*, 2023a.
- 543
- 544 Ziyi Wu, Jingyu Hu, Wuyue Lu, Igor Gilitschenski, and Animesh Garg. SlotDiffusion: Object-  
 545 Centric Generative Modeling with Diffusion Models. *Advances in Neural Information Processing*  
 546 *Systems*, 36:50932–50958, 2023b.
- 547
- 548 Andrii Zadaianchuk, Maximilian Seitzer, and Georg Martius. Object-Centric Learning for Real-  
 549 World Videos by Predicting Temporal Feature Similarities. *Advances in Neural Information Pro-*  
 550 *cessing Systems*, 36, 2024.
- 551
- 552 Rongzhen Zhao, Jian Li, Juho Kannala, and Joni Pajarinen. Predicting video slot attention queries  
 553 from random slot-feature pairs. *arXiv preprint arXiv:2508.22772*, 2025a.
- 554
- 555 Rongzhen Zhao, Vivienne Wang, Juho Kannala, and Joni Pajarinen. Grouped Discrete Representa-  
 556 tion for Object-Centric Learning. In *ECML-PKDD*, 2025b.
- 557
- 558 Rongzhen Zhao, Vivienne Wang, Juho Kannala, and Joni Pajarinen. Multi-Scale Fusion for Object  
 559 Representation. In *ICLR*, 2025c.
- 560
- 561 Rongzhen Zhao, Vivienne Wang, Juho Kannala, and Joni Pajarinen. Vector-Quantized Vision Foun-  
 562 dation Model for Object-Centric Learning. In *ACM Multimedia*, 2025d.
- 563
- 564 Rongzhen Zhao, Yi Zhao, Juho Kannala, and Joni Pajarinen. Slot Attention with Re-Initialization  
 565 and Self-Distillation. In *ACM Multimedia*, 2025e.
- 566

## A APPENDIX

### A.1 LLM USAGE STATEMENT

570 We used GPT-based tools solely for correcting grammar and improving the readability of the  
 571 manuscript. No part of the research ideation, experimental design, analysis, or substantive writ-  
 572 ing was generated by LLMs.

### A.2 ABLATION STUDY

576 We conduct ablation studies as shown in Table 5.

#### (a) Query preheating related:

579 (a.1) Implementing our preheater as a Transformer decoder block is better than as a Slot Attention  
 580 module;

581 (a.1.1) If using a Transformer decoder block as preheater, then switch the self-attention and cross-  
 582 attention in it is better than not;

584 (a.2) Stopping gradient on preheated queries is better than not;

585 (a.3) Setting preheating loss weight to 100 is better than other values;

#### (b) Transform differentiating related:

588 (b.1) Using shared module weights on first-frame transform  $\Phi_a$  and non-first-frame transforms  $\Phi'_a$   
 589 is better than using separate weights;

590 (b.2) For conditioned video OCL, using iteration numbers of 3 and 1 on first and non-first frames  
 591 respectively is better than other combinations;

593 (b.3) For unconditioned video OCL, using iteration numbers of 3 and 1 on first and non-first frames  
 594 respectively is better than other combinations.

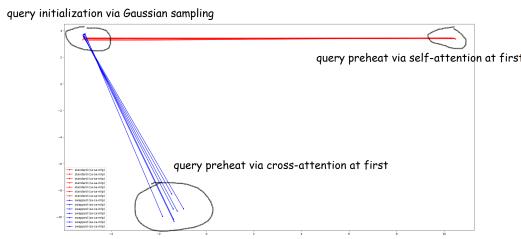
	ARI + ARI <sub>fg</sub>
Preheater implementation @COCO	
a Transformer decoder block	68.3 <sub>±0.8</sub>
a Slot Attention module	63.3 <sub>±1.4</sub>
MLP as the preheater	47.8 <sub>±8.5</sub>
no preheater and preheat loss	56.9 <sub>±2.5</sub>
Switch cross-attention and self-attention in preheater @COCO	
Yes	68.3 <sub>±0.8</sub>
No	49.6 <sub>±9.4</sub>
Stop gradient on preheated query @COCO	
Yes	68.3 <sub>±0.8</sub>
No	67.5 <sub>±2.9</sub>
Preheating loss weight @COCO	
10	59.7 <sub>±1.0</sub>
50	65.5 <sub>±0.4</sub>
100	68.3 <sub>±0.8</sub>
200	67.4 <sub>±1.3</sub>
Use separate weights for first and non-first transforms @YTVIS	
separate	52.3 <sub>±0.7</sub>
shared	68.3 <sub>±0.8</sub>
Unconditional video OCL: first and non-first SA #iter @YTVIS	
3+1	105.6 <sub>±2.2</sub>
1+1	97.4 <sub>±11.4</sub>
3+3	103.4 <sub>±6.8</sub>
Conditional video OCL: first and non-first SA #iter @MOVi-C	
3+1	136.3 <sub>±7.1</sub>
1+1	133.9 <sub>±15.0</sub>
3+3	132.7 <sub>±8.4</sub>

Table 5: Ablation studies.

### A.3 WHY SWAPPING SELF-ATTENTION WITH CROSS-ATTENTION?

Denote self-attention as [sa] while cross-attention as [ca]. The standard Transformer decoder block has the architecture of [sa]-[ca]-[mlp], while our swapped Transformer decoder block has the architecture of [ca]-[sa]-[mlp]. Note that short-cut connections are ignored for simplicity. We use PCA (Principle Component Analysis) to visualize the intermediate slots inside (i) the standard Transformer decoder block, i.e., queries and preheated queries after the [sa] as the first module; and (ii) the swapped Transformer decoder block, i.e., queries and the preheated queries after [ca] as the first module. The model checkpoints are reused from Table 5.

As shown, our swapped attention produces diverse points, i.e., expressive representations.



- the queries after initialization – *query initialization after Gaussian sampling*
  - (top-left): clustered together

- 648 • and the queries after the first attention  
 649   – for standard transformer decoder block as the preheater, the first attention module is  
 650     self attention – *query preheating after self-attention at first*  
 651       \* (top-right): still clustered together  
 652   – for swapped transformer decoder block as the preheater, the first attention module is  
 653     cross attention – *query preheating after cross-attention at first*  
 654       \* (bottom): become well separated  
 655

656 We explain the observed performance gap and visualization as below:  
 657

- 658 • The queries before preheating are sampled randomly from some learnt Gaussian distri-  
 659     butions (the mainstream case), containing no specific information about current specific  
 660     features.  
 661 • Thus using self-attention to mix them is meaningless, as this still does not introduce any  
 662     specific information about current specific visual feature.  
 663 • Thus inside our preheater, we should first inject the current specific information into the  
 664     queries by cross attention, then further transformation like self-attention and MLP could be  
 665     meaningful.  
 666

#### 667 A.4 MATHEMATICAL ANALYSIS

##### 668 Benefit of Preheating

669 Follow the settings and notations from Sections 3.1 and 3.2. Although  $\phi_a$  and  $\phi_p$  take  $\mathbf{F}$  as the  
 670 second input argument, we ignore it for simplicity.

671 In practice, for any fixed  $\mathbf{F}$ ,  $\phi_a$  is usually a contraction, thus for all  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times c}$ ,

$$672 \quad \|\phi_a^i(\mathbf{X}) - \phi_a^i(\mathbf{Y})\| \leq \alpha^i \|\mathbf{X} - \mathbf{Y}\| \quad \text{where } \alpha \in [0, 1] \quad (9)$$

673 By the Banach fixed point theorem, there is a unique fixed point  $\mathbf{S}^* = \phi_a(\mathbf{S}^*)$ , and for every  $\mathbf{X}$ ,

$$674 \quad \|\phi_a^i(\mathbf{X}) - \mathbf{S}^*\| \leq \alpha^i \|\mathbf{X} - \mathbf{S}^*\| \quad (10)$$

675 With our preheater,  $\phi_p(\mathbf{Q})$  is closer to  $\mathbf{S}^*$  than  $\mathbf{Q}$ , then:

$$676 \quad \|\phi_p(\mathbf{Q}) - \mathbf{S}^*\| \leq q \|\mathbf{Q} - \mathbf{S}^*\| \quad \text{where } q \in [0, 1) \quad (11)$$

677 And then:

$$678 \quad \|\phi_a^3(\phi_p(\mathbf{Q})) - \mathbf{S}^*\| \leq \alpha^3 \|\phi_p(\mathbf{Q}) - \mathbf{S}^*\| \leq q \alpha^3 \|\mathbf{Q} - \mathbf{S}^*\| \quad (12)$$

679 Compared with not using preheater,

$$680 \quad \|\phi_a^3(\mathbf{Q}) - \mathbf{S}^*\| \leq \alpha^3 \|\mathbf{Q} - \mathbf{S}^*\| \quad (13)$$

681 Therefore, the preheated run is strictly closer to the fixed point after three iterations than the non-  
 682 preheated run would be.

##### 683 Benefit of Differentiating

684 Follow the settings and notations from Sections 3.1 and 3.3. Although  $\phi_a$  takes  $\mathbf{F}$  as the second  
 685 input argument, we ignore it for simplicity. We treat one frame at a time, dropping subscript  $t$ .

686 We also supplement the decoding part of OCL here, where reconstruction is utilized for supervision.  
 687 Reconstruction is  $\mathbf{X}' = \phi_d(\mathbf{S})$  and loss is  $l = \text{MSE}(\mathbf{X}', \mathbf{X})$ .

688 We follow the assumption of Equation (9).

689 According to Lipschitz Jacobian bounds,

$$690 \quad \left\| \frac{\partial \phi_a(\mathbf{S})}{\partial \mathbf{S}} \right\| \leq \alpha \quad (\text{consistent with contraction}) \quad (14)$$

$$\left\| \frac{\partial \phi_a(\mathbf{S})}{\partial \theta_a} \right\| \leq B \quad \text{for all } \mathbf{S} \text{ (bound on how strong parameters in each iteration)} \quad (15)$$

$$\left\| \frac{\partial \phi_d(\mathbf{S})}{\partial \mathbf{S}} \right\| \leq L \quad (\text{bound on the largest loss}) \quad (16)$$

Unroll the total  $i_1$  iterations. Let  $J_i := \frac{\partial \phi_a(S^{(i-1)})}{\partial S}$  and  $U_i := \frac{\partial \phi_a(S^{(i-1)})}{\partial \theta_a}$ ,  $D := \frac{\partial \phi_d(S)}{\partial S}$  and  $G_X := \frac{\partial l}{\partial X}$ . And the derivative of the final  $S^{(i_1)}$  w.r.t  $\theta_a$  is the sum of contributions from each unrolled iteration:

$$\frac{\mathcal{S}^{(i_1)}}{\partial \theta_a} = \Sigma_{i=1}^{i_1} (\Pi_{m=i+1}^{i_1} \mathbf{J}_m) \mathbf{U}_i \quad (17)$$

By chain rule, the full gradient is

$$\frac{\partial \mathcal{L}}{\partial \theta_a} = \mathbf{G}_{\mathbf{X}}^T \mathbf{D} \frac{\partial \mathbf{S}^{(i_1)}}{\partial \theta_a} = \mathbf{G}_{\mathbf{X}}^T \mathbf{D} \Sigma_{i=1}^{i_1} (\Pi_{m=i+1}^{i_1} \mathbf{J}_m) \mathbf{U}_i \quad (18)$$

For the frame with  $i_1$  iterations,

$$\left\| \frac{\partial l}{\partial \phi_a} \right\| \leq \|G_X\| LB \Sigma_{i=0}^{i_1-1} \alpha^i = \|G_X\| LB \frac{1 - \alpha^{i_1}}{1 - \alpha} < \|G_X\| LB \frac{1}{1 - \alpha} \quad (19)$$

where on the right side, only  $\|G_x\|$  depends on the number of iterations  $i_1$ .

In practice, for the first frame  $\|G_X\|$  tends to be large and more iterations reduces it, while for non-first frames  $\|G_X\|$  tends to be small and less iterations are needed. And if we insist to use more iterations for non-first frames (the same number of iterations as the first frame),  $\|G_X\|$  can be too small compared with that of the first frame. This causes imbalanced gradient contributions from the first and non-first frames to  $\phi_a$  parameters.

## A.5 VIDEO OCL ON YTVIS21 USING ORIGINAL VIDEO LENGTH

Here are the object discovery results on dataset YTVIS21<sup>12</sup>. Note that the results are produced on SlotContrast official codebase<sup>13</sup>, namely, we adopt all the hyperparameters used by SlotContrast, especially using the original video length, instead of clipping them to constant length 20 as in Table 2. Also note that the metrics  $ARI_{f_o}$  and mBO are the video ones, rather than the image ones.

@YTVIS21	video ARIfg	video mBO
VideoSAUR (copied values)	28.9	26.3
SlotContrast (copied values)	38.0	33.7
SmoothSA (seed@42,43,44)	45.9±1.2	36.7±0.5

## A.6 COMPUTATION OVERHEAD

We provide concrete numbers of the computation overhead below. Our method shows better computation efficiency in both training time and memory consumption than baselines.

V100	training time / hours	memory consumption / GB
SPOT @COCO, bs32	4.7	8.5
DIAS @COCO, bs32	4.5	9.4
SmoothSA @COCO, bs32	4.2	8.7
SlotContrast @YTVIS, bs8	7.4	6.4
RandSFQ @YTVIS, bs8	6.8	5.2
SmoothSA @YTVIS, bs8	7.1	5.2

<sup>12</sup><https://youtube-yos.org/challenge/2021>

<sup>13</sup> <https://github.com/martius-lab/slotcontrast>

756  
757

## A.7 RESULTS ON DATASET OVIS

758  
759  
760  
761

We also evaluate our method's effectiveness on more challenging datasets. We choose OVIS<sup>14</sup> (Qi et al., 2022), which is more occluded than YTVIS. Our method still shows some superiority over the baseline. But the performance boosts over the baseline is much smaller, compared with the results of dataset YTVIS in Table 2.

762  
763  
764  
765

@OVIS	ARI	ARI <sub>fg</sub>	mBO	mIoU
SlotContrast	37.2±1.0	44.8±1.0	20.1±0.8	17.1±0.9
SmoothSA	39.3±2.3	47.5±0.4	21.0±0.4	19.5±0.4

766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

---

<sup>14</sup><https://songbai.site/ovis>