

Towards Interpretable Controllability in Object-Centric Learning

Anonymous CVPR submission

Paper ID *****

Abstract

001 *Object-centric learning (OCL) is extensively researched to*
 002 *better understand complex scenes by acquiring object rep-*
 003 *resentations or slots. While recent studies in OCL have*
 004 *made strides with complex images, the interpretability and*
 005 *interactivity over object representation remain largely un-*
 006 *charted. In this paper, we introduce a novel method, Slot*
 007 *Attention with Image Augmentation (SlotAug), to explore*
 008 *the possibility of learning interpretable controllability over*
 009 *slots in a self-supervised manner by utilizing image aug-*
 010 *mentation. We also devise the concept of sustainability in*
 011 *controllable slots by introducing iterative and reversible*
 012 *controls over slots with two proposed submethods: Auxil-*
 013 *iary Identity Manipulation and Slot Consistency Loss.*

014 1. Introduction

015 Compositional comprehension of visual scenes [10, 18, 28],
 016 essential for various computer vision tasks such as local-
 017 ization [4] and reasoning [27], requires human-like under-
 018 standing of complex world [24, 35, 36]. In response to this,
 019 *object-centric learning* (OCL) has emerged as an active re-
 020 search area [11, 22, 25]. OCL aims to enable a model to
 021 decompose an image into objects, and to acquire their rep-
 022 resentations, *slots*, without human-annotated labels.

023 In this work, we advance the field of OCL in terms of
 024 the interpretability of object representation. To achieve in-
 025 terpretable controllability, we propose a method that en-
 026 ables the manipulation of slots through human interpretable
 027 instructions in a self-supervised manner. We address the
 028 training-inference discrepancy by incorporating image aug-
 029 mentation and *slot manipulation* into our training pipeline.
 030 Consequently, we resolve the discrepancy and streamline
 031 the way to interact with slots in the inference.

032 Second, to attain sustainability in object representation,
 033 we introduce *Auxiliary Identity Manipulation* (AIM) and
 034 *Slot Consistency Loss* (SCLoss). AIM is designed to assist
 035 in learning the concept of multi-round manipulation. AIM
 036 is implemented by incorporating an *auxiliary manipulation*
 037 into the intermediate stage of slot manipulation, where the

auxiliary manipulation introduces no semantic changes to
 object properties. This simple auxiliary process can expose
 our model to multi-round manipulation: we can make two-
 round manipulations with one instruction from the augmen-
 tation and the other from the auxiliary manipulation. More-
 over, SCLoss enables learning the concept of reversible ma-
 nipulation, such as the relationship between moving an ob-
 ject to the right and returning to the left. After being trained
 with SCLoss, our model produces consistent and reusable
 representations that can undergo multiple modifications.

Extensive experiments demonstrate the interpretable and
 sustainable controllability of our model. To assess inter-
 pretability, we conduct object manipulation experiments
 where slots are guided by semantically interpretable in-
 structions. In evaluating sustainability, we introduce novel
 experiments, like the durability test. Our evaluation encom-
 passes not only pixel space such as object-level image edit-
 ing, but also slot space such as property prediction, provid-
 ing a comprehensive examination of our method.

2. Method

2.1. Slot Attention with Image Augmentation

Data augmentation. We introduce a simple data aug-
 mentation scheme that, for a given input reference im-
 age $img_{ref} \in \mathbb{R}^{H \times W \times 3}$, generates an augmented image
 $img_{aug} \in \mathbb{R}^{H \times W \times 3}$ and the transformation instructions
 between them, $insts_{ref2aug} \in \mathbb{R}^{K \times L}$, where K and L in-
 dicate the number of slots and the number values for ob-
 ject properties. img_{aug} is produced by randomly translat-
 ing, scaling, or color shifting img_{ref} . To transform img_{ref}
 into img_{aug} , we employ a set of instructions $insts_{ref2aug}$.
 These instructions comprise a list of values that dictate the
 augmentation, including translation values, a scaling factor,
 and color shift values in the HSL color space. We also have
 the inverse instructions, $insts_{aug2ref} \in \mathbb{R}^{K \times L}$, allowing us
 to revert img_{aug} back to img_{ref} . Henceforth, for the sake
 of simplicity in notation, we employ r and a as shorthand
 for ref and aug . For instance, the expression img_r and
 img_{r2a} equal img_{ref} and $img_{ref2aug}$, respectively. More
 details are described in the Appendix.

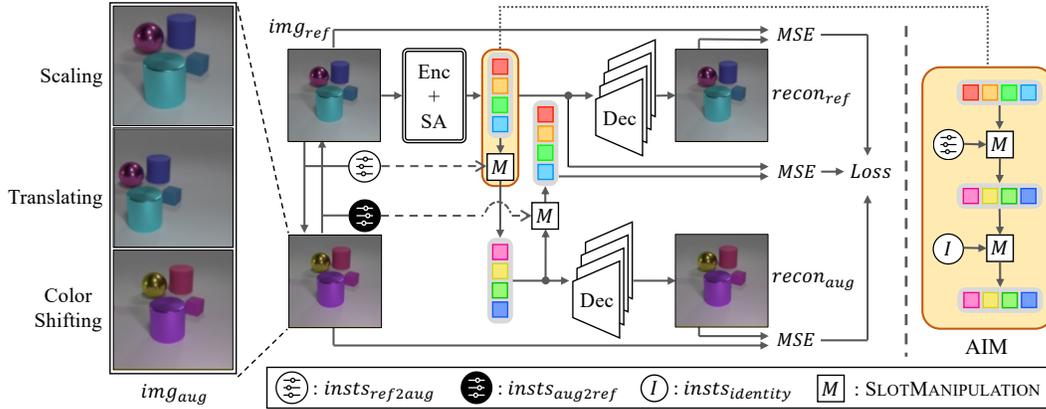


Figure 1. **Model architecture.** From a given image img_{ref} , we generate an augmented image img_{aug} (leftmost part of the figure), and the instruction $insts_{ref2aug}$ and its inverse $insts_{aug2ref}$. Our model produces slots from img_{ref} and decodes them to reconstruct the given image ($recon_{ref}$). The slots are also manipulated, along with $insts_{ref2aug}$, by SLOTMANIPULATION. We incorporate Auxiliary Identity Manipulation (AIM) into this manipulation process. The details are provided in the right part of the figure. The manipulated slots are then simultaneously 1) decoded to the reconstructed augmented image $recon_{aug}$, and 2) re-manipulated by SLOTMANIPULATION with $insts_{aug2ref}$. Our total loss consists of the reconstruction losses of reference and augmented images, and the slot-level cycle consistency.

077 **Training.** We propose a novel training process that leverages
 078 image augmentation (Fig. 1). Our training scheme
 079 enables learning interpretable controllability which allows
 080 us to interact with the model via semantically interpretable
 081 instructions. Our training process involves data augmentation,
 082 spatial binding, slot manipulation, and image reconstruction
 083 via slot decoding. For a given input image, we initially
 084 perform data augmentation to yield img_r , img_a ,
 085 $insts_{r2a}$, and $insts_{a2r}$. Then, the model performs SPA-
 086 TIALBINDING on img_r to produce $slots_{ref}$, or $slots_r$.

087 Thereafter, the model conducts SLOTMANIP (Alg. 1)
 088 to modify $slots_r$ based on $insts_{r2a}$. In SLOTMANIP, we
 089 utilize a newly introduced property encoder denoted by
 090 $PropEnc$ which is 3-layer MLPs. This $PropEnc$ gener-
 091 ates vector representations, $inst_vec$, which capture the
 092 essence of transformation instructions. Each $PropEnc_j$ gener-
 093 ates an $inst_vec_j$ that encodes the values of $insts_{r2a}$
 094 for the j -th property. These vectors are then added to $slots_r$
 095 to reflect the effect of $insts_{r2a}$. This addition is followed by
 096 a residual connection, along with layer normalization and
 097 another MLP to generate $slots_{r2a}$.

098 Lastly, $slots_{r2a}$ is decoded by the decoder to create the
 099 $recon_a$, the reconstruction for the augmented image img_a .
 100 The MSE between img_a and $recon_a$ serves as a training
 101 loss, \mathcal{L}_{aug} . To ensure stable training, we also adopt an addi-
 102 tional loss, \mathcal{L}_{ref} , the MSE between the img_r and $recon_r$, the
 103 reconstructed reference image decoded from $slots_r$. Accord-
 104 ingly, our training loss for image reconstruction is defined
 105 as $\mathcal{L}_{recon} = \mathcal{L}_{ref} + \mathcal{L}_{aug}$.

106 **Inference.** To perform object manipulation, the model
 107 takes the position of the target object, along with the instruc-
 108 tion to be carried out. We use the Hungarian algorithm [23]
 109 to find the slot for the object closest to the given position.

Algorithm 1 Slot manipulation algorithm in pseudo-code. J represents the number of object properties, while $P_{j,f}$ and $P_{j,l}$ indicate the first and last indices of the j -th property values. We use Layer Normalization LN to normalize vectors.

```

1: function SLOTMANIP(slots, insts)
2:   for  $j = 0 \dots J$  do
3:      $inst_j = insts[:, P_{j,f} : P_{j,l}]$ 
4:      $inst\_vec_j = PropEnc_j(LN(inst_j))$ 
5:      $slots = slots + inst\_vec_j$ 
6:   end for
7:    $slots = slots + MLP(LN(slots))$ 
8:   return slots
9: end function

```

To predict the position of an object encoded in a slot, we
 compute the center of mass acquired from the alpha mask
 by the decoder or from the attention map between the visual
 encodings and the slot. After figuring out the desired slot,
 we perform slot manipulation with the given instructions.

2.2. Sustainability in Object Representation

In this work, we introduce *sustainability*: the concept that
 slots should sustain their integrity even after iterative ma-
 nipulations. Thus, sustainability is a key feature that con-
 tributes to the reliable and reusable object representation.

Auxiliary Identity Manipulation (AIM) serves as the
 identity operation for slot manipulation, indicating no
 changes in object properties. By manipulating slots with in-
 structions having zero values for translation, one for scaling,
 and so on, AIM is supposed to make each slot preserve the
 identity of the object. We incorporate AIM into the training
 process to make the model recognize and maintain the in-

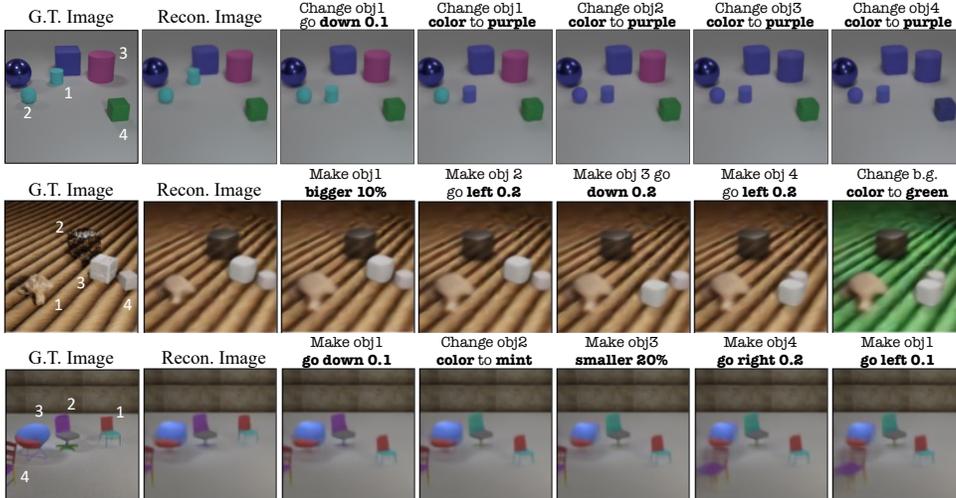


Figure 2. **Results of object manipulation.** The first two columns are the ground truths and reconstructions. The other columns are the results of manipulation along the human-interpretable instructions. The instructions are shown in the text for easy understanding. The instantiation of instruction can be found in the Appendix. From the first row onwards, the results are for CLEVR, CLEVRTEX, and PTR.

intrinsic properties of objects during iterative manipulations. AIM is applied to the slot manipulation process as follows:

$$\begin{aligned} slots'_{r2a} &= f(f(slots_r, insts_{r2a}), insts_{id}) \\ &= f(slots_{r2a}, insts_{id}), \end{aligned} \quad (1)$$

where f represents the SLOTMANIP function, and $insts_{id}$, or $insts_{identity}$, is the instruction that contains the identity elements for manipulating properties. In the followings, $slots'_{r2a}$ is notated as $slots_{r2a}$ for simplicity.

Slot Consistency Loss (SCLoss) addresses the issue of a slot diverging significantly from its original state after iterative manipulations, even when a user intends to restore the corresponding object to its original state. To implement SCLoss, we introduce $slots_{restored}$, which is derived by executing a series of SLOTMANIP operations on $slots_r$ using $insts_{r2a}$ and $insts_{a2r}$. Supposed that our goal is to ensure $slots_r$ and $slots_{restored}$ have the same representation, we set the MSE between them as SCLoss. As a result, the model learns to keep the two distinct slots representing the same object as close as possible and to be robust against multiple rounds of manipulation. The equation of SCLoss, \mathcal{L}_{cycle} , and the total training loss, \mathcal{L}_{total} , are as follows:

$$\mathcal{L}_{cycle} = \frac{1}{K} \|f(slots_{r2a}, insts_{a2r}) - slots_r\|_2^2, \quad (2)$$

$$\mathcal{L}_{total} = w_{recon} \mathcal{L}_{recon} + w_{cycle} \mathcal{L}_{cycle}, \quad (3)$$

where K is the number of slots, f is the SLOTMANIP function, and w_{recon} and w_{cycle} are the weights for the losses.

3. Experiments

Experimental Details. We evaluate models on three multi-object datasets: CLEVR6 [18], CLEVRTEX6 [19] and PTR

[16]. Regarding the baseline models, we basically employ the model architecture of Slot Attention. For CLEVRTEX6 and PTR, we replace the encoder with ViT [5] pretrained by MAE [15] and the decoder with that of SRT [30] while using an increased size of the slot attention module. The additional details for adopting large models are described in the Appendix. To clarify the methods used in ablation studies, we categorize our model into three versions: **v1**, which is exclusively trained with image augmentation; **v2**, which improves upon v1 with AIM; **v3**, which extends v1 with both AIM and SCLoss. More details including the training schemes are described in the Appendix.

3.1. Interpretable Controllability: Image Editing

As shown in Fig. 2, we can manipulate various properties of a single object, even multiple times, according to user-given instructions. Based on this observation, we can ascertain that our object representations retain the intrinsic properties of objects seamlessly even after manipulation. This interpretable controllability is achieved with a neglectable compromise of the performance on the object discovery task. It is also worth noting that we can also manipulate the background as shown in the second row of Fig. 2 since we employ SLASH [20] which treats the background as a single entity. Deeper discussions with theoretical proof and empirical results are provided in the Appendix.

3.2. Sustainability: Durability Test

To assess the sustainability of our model, we devise a novel evaluation, called *durability test*. In the durability test, we evaluate how many manipulations a model can endure while preserving object representation intact. As shown in Fig.

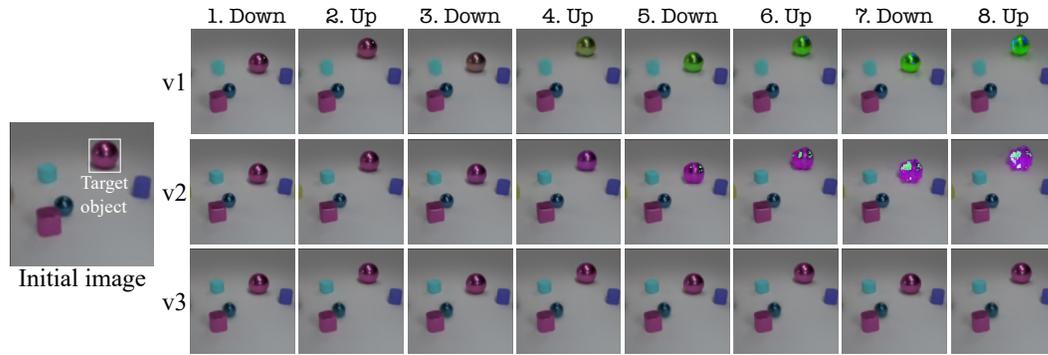


Figure 3. **Durability test.** The leftmost image is the initial image for the test. Each row shows the results of that each model is instructed to alternately move the target object up and down four times each.

Table 1. **Results on durability test** with MSE on CLEVR6.

	Slot (\downarrow)	Obj. Pos. (\downarrow)
	Single step (x8)	
(v1) Train w/ aug.	50.8	0.14
(v2) + AIM	39.7	0.15
(v3) + AIM + SCLoss	0.25	0.01
	Multiple steps (x4)	
(v1) Train w/ aug.	54.0	0.16
(v2) + AIM	41.4	0.11
(v3) + AIM + SCLoss	0.31	0.02

185 3, while v1 fails to keep the color after the second round,
 186 v2 relatively preserves the color well for the fourth round.
 187 Nevertheless, from the fifth round, the texture progressively
 188 diverges from its original. Different from the v1 and v2, v3
 189 demonstrates strong durability despite a greater number of
 190 manipulations. We also perform quantitative evaluations on
 191 100 randomly selected samples from CLEVR6 to measure
 192 the intrinsic deformity of slots and extrinsic change of
 193 object properties, especially position. As shown in Tab. 1, we
 194 achieve better sustainability as the model evolves from v1
 195 to v2 and v3 by utilizing the proposed AIM and SCLoss.

196 3.3. Slot Space Analysis: Property Prediction

197 In addition to the pixel space analysis, for a comprehensive
 198 assessment of the effectiveness of our method, we extend our
 199 examination to the analysis of the latent slot space. To evaluate
 200 the quality of slots concerning human-interpretable object
 201 properties, such as size, color, material, shape, and position,
 202 we conduct a *property prediction* task using CLEVR6. This
 203 task enables us to scrutinize how well the slots are distributed
 204 within the latent vector space following the properties of
 205 corresponding objects. The analysis allows for an understanding
 206 of object representations that impart semantic existence
 207 significance beyond mere segmenting objects in an image.
 208

209 A *property predictor*, consisting of 3-layer MLPs, takes
 210 slots as input and predicts a property of objects. Each prop-

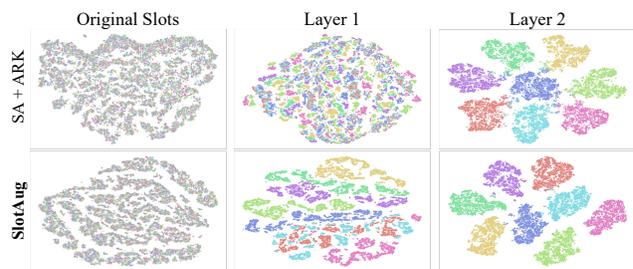


Figure 4. **t-SNE of slots on property prediction for color.** Each row shows the results of the baseline model (SA + ARK) and our model (SlotAug), respectively. The first column is the result of the original slots obtained from the spatial binding. The second and third columns are the results of the intermediate outputs from the first and second MLP layers of the property predictor, respectively.

Table 2. **Property prediction results** on CLEVR6 in F1 score. The number of classes are shown in the parenthesis after the property name. We set two distance thresholds for position prediction.

	SA + ARK	SlotAug (Ours)
Size (2)	69.7	82.2
Color (8)	63.5	78.2
Material (2)	70.4	82.6
Shape (3)	59.1	73.0
Pos@0.15	71.1	84.2
Pos@0.05	51.8	77.2

211 erty predictor is trained by supervised learning using the
 212 ground truths matched by the Hungarian algorithm [23].
 213 To assess the quality of object representations, we freeze
 214 the OCL models that produce slots. As shown in Tab. 2,
 215 our model outperforms the baseline method [20] across all
 216 properties including those, like material and shape, that are
 217 not addressed during training. Moreover, in Fig. 4, qualitative
 218 results using t-SNE [37] show that while the original
 219 slots do not appear to be well-clustered, slots obtained by
 220 SlotAug exhibit better adaptability to the downstream task.
 221 This observation reinforces our quantitative findings.

222

References

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

- [1] Pierre Baldi. Autoencoders, unsupervised learning, and deep architectures. In *Proceedings of ICML workshop on unsupervised and transfer learning*, pages 37–49. JMLR Workshop and Conference Proceedings, 2012. 6
- [2] Ondrej Biza, Sjoerd van Steenkiste, Mehdi SM Sajjadi, Gamaleldin F Elsayed, Aravindh Mahendran, and Thomas Kipf. Invariant slot attention: Object discovery with slot-centric reference frames. *arXiv preprint arXiv:2302.04973*, 2023. 6, 8
- [3] Christopher P Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt Botvinick, and Alexander Lerchner. Monet: Unsupervised scene decomposition and representation. *arXiv preprint arXiv:1901.11390*, 2019. 6
- [4] Minsu Cho, Suha Kwak, Cordelia Schmid, and Jean Ponce. Unsupervised object discovery and localization in the wild: Part-based matching with bottom-up region proposals. In *CVPR*, pages 1201–1210, 2015. 1
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. 3, 8
- [6] Gamaleldin Elsayed, Aravindh Mahendran, Sjoerd van Steenkiste, Klaus Greff, Michael C Mozer, and Thomas Kipf. Savi++: Towards end-to-end object-centric learning from real-world videos. *Advances in Neural Information Processing Systems*, 35:28940–28954, 2022. 8
- [7] Martin Engelcke, Adam R Kosiorek, Oiwi Parker Jones, and Ingmar Posner. Genesis: Generative scene inference and sampling with object-centric latent representations. *arXiv preprint arXiv:1907.13052*, 2019. 6
- [8] Martin Engelcke, Oiwi Parker Jones, and Ingmar Posner. Genesis-v2: Inferring unordered object representations without iterative refinement. *Advances in Neural Information Processing Systems*, 34:8085–8094, 2021. 6
- [9] Jerome Feldman. The neural binding problem (s). *Cognitive neurodynamics*, 7:1–11, 2013. 6
- [10] Martin A Fischler and Robert A Elschlager. The representation and matching of pictorial structures. *IEEE Transactions on computers*, 100(1):67–92, 1973. 1
- [11] Klaus Greff, Antti Rasmus, Mathias Berglund, Tele Hao, Harri Valpola, and Jürgen Schmidhuber. Tagger: Deep unsupervised perceptual grouping. *Advances in Neural Information Processing Systems*, 29, 2016. 1
- [12] Klaus Greff, Raphaël Lopez Kaufman, Rishabh Kabra, Nick Watters, Christopher Burgess, Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-object representation learning with iterative variational inference. pages 2424–2433. PMLR, 2019. 6
- [13] Klaus Greff, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. On the binding problem in artificial neural networks. *arXiv preprint arXiv:2012.05208*, 2020. 6
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceed-*

ings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016. 8 279
280

- [15] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022. 3, 8 281
282
283
- [16] Yining Hong, Li Yi, Joshua B Tenenbaum, Antonio Torralba, and Chuang Gan. Ptr: A benchmark for part-based conceptual, relational, and physical reasoning. *NeurIPS*, 2021. 3 284
285
286
- [17] Daniel Im Im, Sungjin Ahn, Roland Memisevic, and Yoshua Bengio. Denoising criterion for variational auto-encoding framework. In *Proceedings of the AAAI conference on artificial intelligence*, 2017. 9 287
288
289
290
- [18] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, pages 2901–2910, 2017. 1, 3 291
292
293
294
295
- [19] Laurynas Karazija, Iro Laina, and Christian Rupprecht. Clevrtex: A texture-rich benchmark for unsupervised multi-object segmentation. *NeurIPS*, 2021. 3 296
297
298
- [20] Jinwoo Kim, Janghyuk Choi, Ho-Jin Choi, and Seon Joo Kim. Shepherding slots to objects: Towards stable and robust object-centric learning. *arXiv preprint arXiv:2303.17842*, 2023. 3, 4, 8, 9 299
300
301
302
- [21] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 6 303
304
- [22] Thomas Kipf, Gamaleldin F Elsayed, Aravindh Mahendran, Austin Stone, Sara Sabour, Georg Heigold, Rico Jonschkowski, Alexey Dosovitskiy, and Klaus Greff. Conditional object-centric learning from video. *arXiv preprint arXiv:2111.12594*, 2021. 1, 6 305
306
307
308
309
- [23] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 1955. 2, 4 310
311
- [24] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *Behavioral and brain sciences*, 40:e253, 2017. 1 312
313
314
315
- [25] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. 33:11525–11538, 2020. 1, 6, 8 316
317
318
319
320
- [26] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015. 6 321
322
323
- [27] Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B Tenenbaum, and Jiajun Wu. The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision. *arXiv preprint arXiv:1904.12584*, 2019. 1 324
325
326
327
- [28] David Marr. *Vision: A computational investigation into the human representation and processing of visual information*. MIT press, 2010. 1 328
329
330
- [29] Mehdi SM Sajjadi, Daniel Duckworth, Aravindh Mahendran, Sjoerd van Steenkiste, Filip Pavetic, Mario Lucic, Leonidas J Guibas, Klaus Greff, and Thomas Kipf. Object scene representation transformer. 35:9512–9524, 2022. 6, 8 331
332
333
334

- 335 [30] Mehdi SM Sajjadi, Henning Meyer, Etienne Pot, Urs
336 Bergmann, Klaus Greff, Noha Radwan, Suhani Vora, Mario
337 Lučić, Daniel Duckworth, Alexey Dosovitskiy, et al. Scene
338 representation transformer: Geometry-free novel view syn-
339 thesis through set-latent scene representations. In *Proceed-*
340 *ings of the IEEE/CVF Conference on Computer Vision and*
341 *Pattern Recognition*, pages 6229–6238, 2022. 3, 8
- 342 [31] Maximilian Seitzer, Max Horn, Andrii Zadaianchuk, Do-
343 minik Zietlow, Tianjun Xiao, Carl-Johann Simon-Gabriel,
344 Tong He, Zheng Zhang, Bernhard Schölkopf, Thomas Brox,
345 et al. Bridging the gap to real-world object-centric learning.
346 *arXiv preprint arXiv:2209.14860*, 2022. 6
- 347 [32] Gautam Singh, Fei Deng, and Sungjin Ahn. Illiterate dall-e
348 learns to compose. *arXiv preprint arXiv:2110.11405*, 2021.
349 6, 9, 11
- 350 [33] Gautam Singh, Yi-Fu Wu, and Sungjin Ahn. Simple unsu-
351 pervised object-centric learning for complex and naturalistic
352 videos. *Advances in Neural Information Processing Systems*,
353 35:18181–18196, 2022. 6
- 354 [34] Gautam Singh, Yeongbin Kim, and Sungjin Ahn. Neural
355 systematic binder. In *ICLR*, 2023. 6
- 356 [35] Elizabeth S Spelke and Katherine D Kinzler. Core knowl-
357 edge. *Developmental science*, 10(1):89–96, 2007. 1
- 358 [36] Anne Treisman. The binding problem. *Current opinion in*
359 *neurobiology*, 6(2):171–178, 1996. 1, 6
- 360 [37] Laurens Van der Maaten and Geoffrey Hinton. Visualizing
361 data using t-sne. *Journal of machine learning research*, 9
362 (11), 2008. 4
- 363 [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko-
364 reit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia
365 Polosukhin. Attention is all you need. 30, 2017. 9
- 366 [39] Yanbo Wang, Letao Liu, and Justin Dauwels. Slot-vae:
367 Object-centric scene generation with slot attention. 2023.
368 6
- 369 [40] Nicholas Watters, Loic Matthey, Christopher P Burgess, and
370 Alexander Lerchner. Spatial broadcast decoder: A simple ar-
371 chitecture for learning disentangled representations in vaes.
372 *arXiv preprint arXiv:1901.07017*, 2019. 8
- 373 [41] Ziyi Wu, Nikita Dvornik, Klaus Greff, Thomas Kipf, and
374 Animesh Garg. Slotformer: Unsupervised visual dynam-
375 ics simulation with object-centric models. *arXiv preprint*
376 *arXiv:2210.05861*, 2022. 6
- 377 [42] Sirui Xie, Ari S Morcos, Song-Chun Zhu, and Ramakr-
378 ishna Vedantam. Coat: Measuring object compositionality
379 in emergent representations. pages 24388–24413. PMLR,
380 2022. 6

A. Related Works

The binding problem in artificial neural networks [13], in-
spired by cognitive science [9, 36], is a subject of active
exploration, aiming to attain human-like recognition abil-
ities by understanding the world in terms of symbol-like
entities such as objects. In computer vision, object-centric
learning (OCL) focuses on comprehending visual scenes by
considering objects and their relationships without labeled
data [8, 41, 42]. MONet [3], IODINE [12], and GENE-
SIS [7] have adopted autoencoding architectures [1, 21, 26]
to accomplish self-supervised OCL, and Slot Attention [25]
introduced the concept of slot competition, which enables
parallel updates of slots with a single visual encoding and
decoding stage. Recent studies have leveraged large-scale
models to learn object representations in complex images
[31, 32], multi-view images [29], and videos [22, 33].

Several studies have shown the possibility of interact-
ing with object representation to manipulate the objects.
VAE-based models such as IODINE [12] and Slot-VAE [39]
showed that adjusting the values of slots can change ob-
ject properties. SysBinder [34] demonstrated that replacing
factor-level slot, called block, between slots exchanges the
corresponding properties. However, these works have diffi-
culties in determining ways to interact with slots as they re-
quire manual efforts to identify the features associated with
specific properties. ISA [2] incorporates spatial symmetries
of objects using slot-centric reference frames into the spa-
tial binding process, enhancing interactivity of object rep-
resentation for spatial properties such as position and scale.
Meanwhile, our method itself has no constraint on the types
of the target property, showing its expandability toward ex-
trinsic properties such as the shape and material of objects
if there exist proper image augmentation skills or labeled
data.

B. Spatial Binding in Slot Attention

The core mechanism of the slot attention, the spatial bind-
ing, is described in Alg. 2. Given an input image img
 $\in \mathbb{R}^{H \times W \times 3}$, CNN encoder generates a visual feature map
 $\text{input} \in \mathbb{R}^{N \times D_{\text{enc}}}$, where H , W , N , and D_{enc} are the
height and width of the input image, the number of pixels
in the input image ($= HW$), and the channel of the visual
feature map. The slot attention module takes slots and
 inputs , and projects them to dimension D_{slot} through
linear transformations k , q , and v . Dot-product attention
is applied to generate an attention map, attn , with query-
wise normalized coefficients, enabling slots to compete for
the most relevant pixels of the visual feature map. The at-
tention map coefficients weight the projected visual feature
map to produce updated slots, updates . With the itera-
tive mechanism of the slot attention module, the slots can
gradually refine their representations.

Algorithm 2 Spatial binding in slot attention algorithm in pseudo-code format. The input image is encoded into a set of N vectors of dimension D_{input} which is mapped to a set of K vectors with dimension D_{slot} . Slots are initialized from a Gaussian distribution with learned parameters $\mu, \sigma \in \mathbb{R}^{D_{slot}}$. The number of iterations is set to $T = 3$.

```

1: function SPATIALBINDING( $img \in \mathbb{R}^{H \times W \times 3}$ )
2:   inputs = Encoder(img)
3:   inputs = LayerNorm(inputs)
4:   for  $t = 0 \dots T$  do
5:     slots_prev = slots
6:     slots = LayerNorm(slots)
7:     attn = Softmax( $\frac{1}{\sqrt{D_{slot}}} k(inputs) \cdot q(slots)^T$ , axis='slots')
8:     updates = WeightedMean(weights=attn $\epsilon$ , values= $v(inputs)$ )
9:     slots = GRU(state=slots_prev, inputs=updates)
10:    slots = slots + MLP(LayerNorm(slots))
11:  end for
12:  return slots
13: end function

```

432 C. Implementation and experimental details

433 C.1. Training

434 We use a single V100 GPU with 16GB of RAM with 1000
 435 epochs and a batch size of 64. The training takes approxi-
 436 mately 65 hours (wall-clock time) using 12GB of RAM for
 437 the CLEVR6 dataset, and 22 hours using 9GB of RAM for
 438 the Tetrominoes dataset, both with 16-bit precision.

439 C.2. Image Augmentation

440 Upon receiving an input image img_{input} , we produce four
 441 outputs: a reference image, denoted as img_{ref} , an aug-
 442 mented image, represented as img_{aug} , and the transforma-
 443 tion instructions between them, indicated as $insts_{ref2aug}$
 444 and $insts_{aug2ref}$.

445 In the data augmentation process, three pivotal variables
 446 are defined. The first is the template size \mathcal{T} , employed for
 447 the initial cropping of img_{input} prior to the application of
 448 transformation (240 for CLEVR6 and 80 for Tetrominoes).
 449 Next, the crop size \mathcal{C} is used to crop the transformed image
 450 before resizing it to \mathcal{M} (192 for CLEVR6 and 64 for Tetro-
 451 minoes). This two-stage cropping procedure mitigates the
 452 zero-padding that results from transformations. Lastly, the
 453 image size \mathcal{M} denotes the final image size post data aug-
 454 mentation (128 for CLEVR6 and 64 for Tetrominoes).

455 In the training phase, img_{ref} is obtained by applying a
 456 center-crop operation on img_{input} using \mathcal{C} and then resiz-
 457 ing it to \mathcal{M} . The generation of img_{aug} is more complex,
 458 entailing the application of a random transformation from a
 459 set of three potential transformations. Initially, img_{input} is
 460 cropped using \mathcal{T} , and the transformation process is imple-
 461 mented. Following this, the transformed image is cropped
 462 by \mathcal{C} and then resized to \mathcal{M} , yielding img_{aug} . The detailed
 463 description for each transformation is as follows:

464 **Translating.** We set a maximum translation value $d_{max} =$

$\frac{\mathcal{T}-\mathcal{C}}{2}$. A value is randomly chosen within the range of
 ($-d_{max}, d_{max}$) for translation along the x -axis (d_x) and the
 y -axis (d_y) respectively.

Scaling. The maximum and minimum scaling factors, s_{max}
 and s_{min} , are computed by $\frac{\mathcal{T}}{\mathcal{C}}$ and $\frac{\mathcal{C}}{\mathcal{T}}$, respectively. A float
 value s , serving as a scaling factor, is then randomly sam-
 pled from within the range of (s_{max}, s_{min}) . One thing to
 note is that calculating the transformation instructions is
 not straightforward due to the potential translation of ob-
 jects during scaling. Thus, to calibrate the instructions, we
 infer translation values from the predicted object positions
 before scaling. The position prediction is calculated as the
 weighted mean on the attention maps between the visual en-
 codings and slots. With this position prediction, we add the
 translation term into the scaling process so that the model
 should perform both object-level scaling and translating:
 $\vec{d} = (s - 1)(\vec{p} - \vec{c})$, where \vec{d} represents the vector of the
 translation value, \vec{p} refers to the vector of the predicted ob-
 ject position, and \vec{c} is the vector corresponding to the posi-
 tion of image center.

Color shifting. In this study, we employ the HSL (hue, sat-
 uration, and lightness) color space for effective object color
 manipulation. The input image, initially in RGB space, is
 converted to HSL space. We adjust the hue by rotating it us-
 ing randomly sampled angles that span the entire hue space.
 For saturation, we apply a scaling factor, determined by the
 exponential of a value randomly drawn from $(-1, 1)$, a hyper
 parameter. Our primary focus lies on the internal color of
 objects, leaving lightness untouched. Nonetheless, adjust-
 ments to lightness can be made if necessary.

Instruction. Each transformation instruction is a list of 6
 values: one scaling factor (Λ_{scale}), two translation param-
 eters ($\Delta x, \Delta y$), and three color shifting parameters in HSL
 ($\Delta_{hue}, \Lambda_{saturation}, \Lambda_{lightness}$) where Λ and Δ means the
 multiplicative and additive factor for the corresponding val-

ues, respectively. The identity instruction, $inst_{identity}$, contains the base values for each transformation. Thus, $inst_{identity}$ has 1 for scaling, (0,0) for translation, and (0,1,1) for color shifting. For the inverse instruction, $inst_{aug2ref}$ has the values of $-inst_{ref2aug}$ for additive factors, and $\frac{1}{inst_{ref2aug}}$ for multiplicative factors.

506 C.3. Model

507 Basically, our model framework is built on Slot Attention
508 [25], thereby the encoder, decoder, and slot attention mod-
509 ule are the same as that of Slot Attention except for the
510 inclusion of the Attention Refining Kernel (ARK) from
511 SLASH [20]. For Tetriminos and CLEVR, we employ a
512 4-layer CNN encoder and a 6-layer Spatial Broadcast (SB)
513 decoder [40] with a hidden dimension of 64. Within the slot
514 attention module, we set the slot dimension to 64, perform
515 the binding process for 3 iterations, and use a kernel size of
516 5 for the ARK. Please refer to the original papers [20, 25]
517 for additional details for Slot Attention.

518 For CLEVRTEX6 and PTR datasets which include more
519 complicated objects, we adopt larger models with a slot di-
520 mension, D_{slot} , of 256. As encoders, we use 1) Resnet34
521 [14] following [2, 6] and 2) ViT-base [5], with the patch
522 size of 8, pretrained via MAE [15] As decoders, we use an
523 increased size of SB decoder consisting of 8-layer CNNs
524 with a hidden dimension of 128, and a Transformer-based
525 decoder proposed in SRT [30]. The original SRT decoder
526 is designed to operate at the image level, and the follow-
527 ing research OSRT [29] introduce a modification to decode
528 slots simultaneously. In this paper, we slightly modified it to
529 decode each slot independently following the spatial broad-
530 cast decoder. This selection is made to demonstrate that
531 our proposed method is not limited to CNN-based spatial
532 broadcast decoders used in Slot Attention but can robustly
533 operate within transformer-based decoders as well, given
534 the appropriate conditions for independence.

535 In Alg. 1 of the main paper, the Property Encoder
536 (PropertyEncoder) takes as input the values that corre-
537 spond to specific properties. Accordingly, the input size for
538 the property encoder is 1 for scaling, 2 for translation, and
539 3 for color shifting. Each property is encoded via Property
540 Encoder, a 3-layer MLP with ReLU activation functions,
541 resulting in a $inst_vec$, a vector of dimension D_{slot} .

542 D. Discussion: How does it work?

543 To begin with, we would like to highlight our unique ap-
544 proach to the training procedure. While our training incor-
545 porates manipulations at the *image-level*, it can be perceived
546 as training the model at the individual *object-level*. In this
547 section, we discuss on how this transition is achieved with-
548 out the need for an additional tuning process, and present
549 empirical results that support our claim.

As we discussed shortly in Sec. 3.1. in the main paper,
the success of transitioning from image-level augmentation
during training to object-level manipulation during infer-
ence can be attributed primarily to the fact that the entire
process for each slot, including object discovery and de-
coding, exclusively influences the reconstruction of its re-
spective *object*. A mathematical proof is provided below to
show how an image-level reconstruction loss can be disen-
tangled into object-level reconstruction losses.

$$\mathcal{L}_{recon} = \|\hat{\mathcal{I}} - \mathcal{I}\|_2^2 \quad (4) \quad 559$$

$$= \left\| \sum_{k=1}^{\mathcal{K}} (\hat{\mathcal{I}}_k^{rgb} \odot \hat{\mathcal{I}}_k^\alpha) - \mathcal{I} \right\|_2^2 \quad (5) \quad 560$$

$$= \left\| \sum_{k=1}^{\mathcal{K}} (\hat{\mathcal{I}}_k^{rgb} \odot \hat{\mathcal{I}}_k^\alpha) - \sum_{k=1}^{\mathcal{K}} (\mathcal{I} \odot \hat{\mathcal{I}}_k^\alpha) \right\|_2^2 \quad (6) \quad 561$$

$$= \left\| \sum_{k=1}^{\mathcal{K}} (\hat{\mathcal{I}}_k^{rgb} \odot \hat{\mathcal{I}}_k^\alpha - \mathcal{I} \odot \hat{\mathcal{I}}_k^\alpha) \right\|_2^2 \quad (7) \quad 562$$

$$\approx \left\| \sum_{k=1}^{\mathcal{K}} (\hat{\mathcal{O}}_k - \mathcal{O}_k) \right\|_2^2, \quad (8) \quad 563$$

$$= \sum_{k=1}^{\mathcal{K}} \left\| (\hat{\mathcal{O}}_k - \mathcal{O}_k) \right\|_2^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^{\mathcal{K}} (\hat{\mathcal{O}}_i \cdot \hat{\mathcal{O}}_j - 2 \hat{\mathcal{O}}_i \cdot \mathcal{O}_j + \mathcal{O}_i \cdot \mathcal{O}_j) \quad (9) \quad 564$$

$$\approx \sum_{k=1}^{\mathcal{K}} \left\| (\hat{\mathcal{O}}_k - \mathcal{O}_k) \right\|_2^2, \quad (10) \quad 565$$

where \mathcal{K} is the number of slots, $\hat{\mathcal{I}} \in \mathbb{R}^{H \times W \times 3}$ represents
the reconstructed image, and $\mathcal{I} \in \mathbb{R}^{H \times W \times 3}$ represents the
input image. $\hat{\mathcal{I}}_k^{rgb} \in \mathbb{R}^{H \times W \times 3}$ and $\hat{\mathcal{I}}_k^\alpha \in \mathbb{R}^{H \times W \times 1}$ are the
reconstruction results generated by the decoder using the k -
th slot as input: an RGB and an alpha map (or an attention
mask), respectively. $\hat{\mathcal{O}}_k \in \mathbb{R}^{H \times W \times 3}$ is the predicted image
for the specific object that is bounded with the k -th slot,
while $\mathcal{O}_k \in \mathbb{R}^{H \times W \times 3}$ is the corresponding ground-truth
object image.

From Eq. (4) to Eq. (5), we follow the decoding pro-
cess of Slot Attention [25]. In particular, each k -th slot is
decoded independently, resulting in the reconstructed RGB
image $\hat{\mathcal{I}}_k^{rgb}$ and the reconstructed alpha map $\hat{\mathcal{I}}_k^\alpha$. The final
reconstruction image $\hat{\mathcal{I}}$ is generated by aggregating $\hat{\mathcal{I}}_k^{rgb}$ us-
ing a pixel-level weighted average, where the weights are
determined by $\hat{\mathcal{I}}_k^\alpha$. It is crucial to recognize that $\hat{\mathcal{I}}_k^\alpha$ serves
as an attention mask, as elaborated below:

$$\sum_{k=1}^{\mathcal{K}} \hat{\mathcal{I}}_k^\alpha(x, y) = 1 \quad \text{for all } x, y, \quad (11) \quad 583$$

where $\hat{\mathcal{I}}_k^\alpha(x, y)$ is a value for the position (x, y) . This characteristic plays a pivotal role in our approach, facilitating the transition from Eq. (5) to Eq. (6). In this transformation, the input image \mathcal{I} is effectively weighted by the set of \mathcal{K} alpha maps, denoted as $\hat{\mathcal{I}}_k^\alpha$, where k spans from 1 to \mathcal{K} . Then, as both the first and second terms in Eq. (6) involve the same sigma operations, we can simplify the expression by combining the individual subtraction operations into a single sigma operation (Eq. (7)).

Subsequently, we approximate Eq. (7) as Eq. (8) to get an object-level disentangled version of the reconstruction loss. Here we assume that both $\hat{\mathcal{O}}_k$ and \mathcal{O}_k only consist of a specific region of interest within the input image. This region corresponds to the target object which is bound to the k -th slot, while the remaining areas are masked out and assigned a value of zero. We can make this assumption based on the successful performance of the previous object-centric learning model, SLASH [20]. SLASH has demonstrated effective capabilities in focusing on and capturing specific objects of interest within an image, by introducing the Attention Refining Kernel (ARK). By incorporating ARK into our model, we confidently assume that $\hat{\mathcal{O}}_k$ and \mathcal{O}_k primarily represent the target object while masking out other irrelevant parts as zero as shown in Fig. 5.

Here, we would like to note that ARK is an optional component in our method, not a necessity. The use of ARK is not intended to enhance object discovery performance in a single training session; rather, it is employed to ensure consistent results across multiple experiments. If our proposed training scenario arises where bleeding issues do not occur in the original SA, it can be achieved without the need for ARK. To substantiate this claim, we present qualitative results in Fig. 7, where we train SlotAug with the original SA (without ARK). One can easily catch that the object manipulation fails in the case of bleeding problem. Specifically, the analysis for the failure case in bleeding problem is as follows: 1) Obviously, if the attention map corresponding to the target object encompasses other objects, it becomes impossible to exclusively manipulate solely the target object, leading to unexpected artifacts in other objects. 2) Whenever tinting instructions are applied, objects become gray and we attribute this to the backgrounds – having a gray color – intervening with the target objects during training.

Eq. (8) can be broken down into two separate summations. The first one is our target term that is the sum of object-level MSE losses, and the second term is the residual term. Lastly, the transition from Eq. (9) to Eq. (10) constitutes a significant simplification in the representation of the loss function. This is a valid transformation under the assumption follows:

$$\hat{\mathcal{O}}_i \cdot \hat{\mathcal{O}}_j = \hat{\mathcal{O}}_i \cdot \mathcal{O}_j = \mathcal{O}_i \cdot \mathcal{O}_j = 0 \quad \text{if } i \neq j. \quad (12)$$

This assumption postulates that the inner product of dif-

ferent object images, whether they are predicted or ground-truth, is always zero. We assert that this assumption is justifiable, much like the previous one, given the promising results obtained in our object discovery experiments. The loss computation is thus decomposed into individual components for each slot, which lends itself to an interpretation of object-level loss.

The conversion from image-level MSE loss to a sum of individual object-level MSE losses provides a new perspective on our training method. Despite the use of image-level manipulations, the underlying core of the training process inherently engages with object-level representations. This demonstrates how a simple methodological addition, incorporating image augmentation into the training process, can lead to considerable gains in the model’s capacity for user-intention-based object manipulation.

Fig. 5 empirically demonstrates the effectiveness of our model, leveraging Slot Attention for controllability over slots. Conversely, it was noted that the well-known alternative framework for object-centric learning, SLATE [32], employing image tokenization from Discrete VAE (dVAE) [17] and Transformer-based auto-regressive decoding [38], struggled with the manipulation of slots, as illustrated in Fig. 6. The same slot manipulation strategy via Property Encoder was used for comparison. Other training environments are just the same as the official paper [32] except for the addition of the training loss for the reconstruction of the augmented images.

E. Conclusion

We presented an OCL framework, SlotAug, for exploring the potential of interpretable controllability in object-centric learning. To achieve this goal, we tackled the object manipulation task, where we added some conditions regarding interpretability and interactivity, via controlling object representations called slots. We employed image augmentation for training our model in a self-supervised manner to resolve the lack of labeled data. Moreover, we introduced a concept of sustainability in slots, achieved by the proposed method AIM and SCLoss. We substantiated the effectiveness of our methods by providing extensive empirical studies and theoretical evidence in the Appendix. These empirical studies include pixel- and slot-space analyses on tasks such as the durability test and property prediction. Though our work remains several questions detailed in the Appendix and represents just one step on a long journey of OCL, we firmly believe that our work is a foundational piece in the field of interpretable OCL and propel the ongoing effort to equip machines with human-like comprehension abilities.

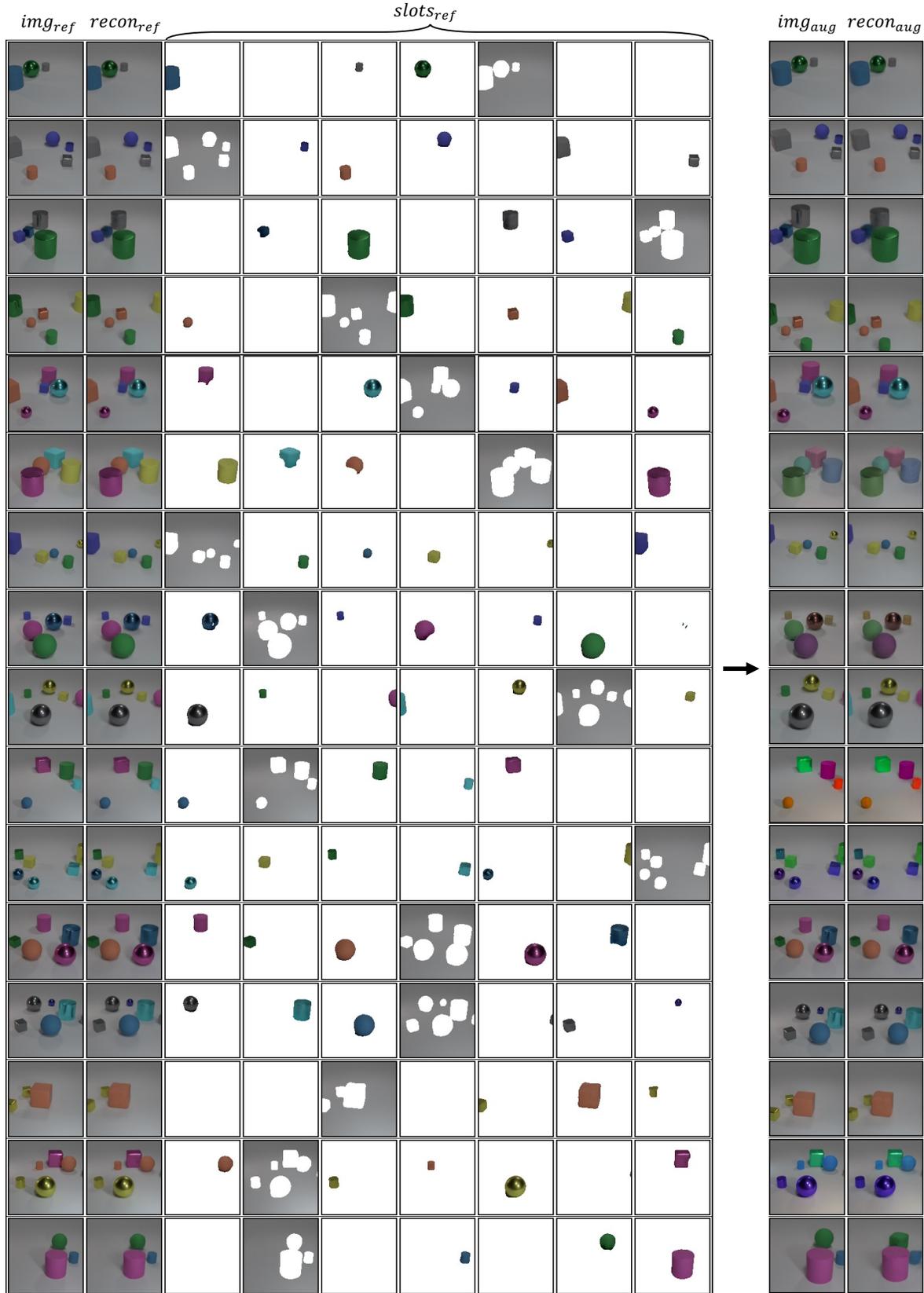


Figure 5. **Training results of our method.** The leftmost column is the reference images, img_{ref} . The second leftmost column is the reconstruction of the reference images, $recon_{ref}$. The middle columns show the object discovery results where each column corresponds to a single slot in $slots_{ref}$. The second rightmost column is the augmented images, img_{aug} . The rightmost column is the reconstruction of the augmented images, $recon_{aug}$.

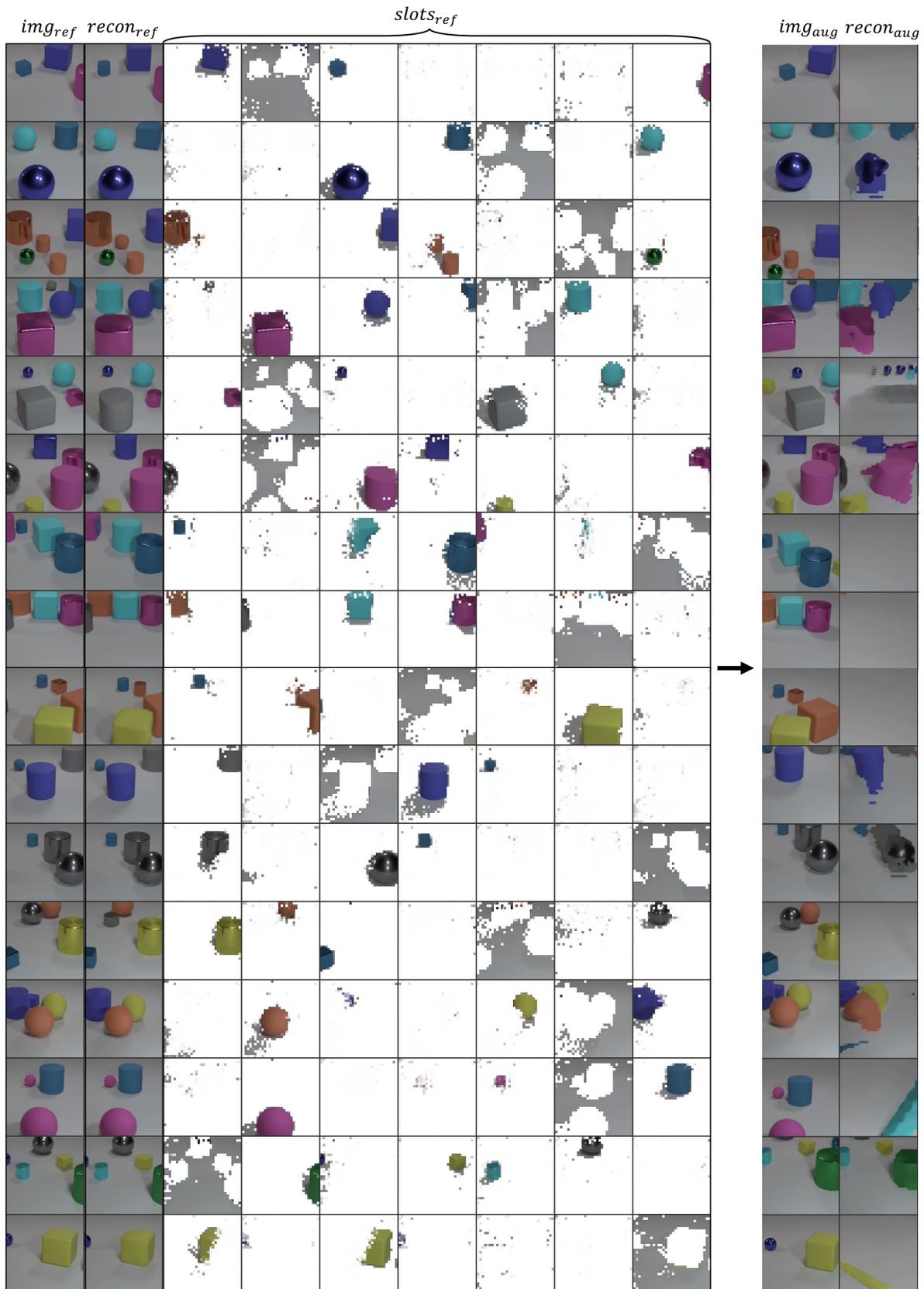


Figure 6. **Training results of SLATE [32] for slot manipulation.** The leftmost column is the reference images, img_{ref} ; The second leftmost column is the reconstruction of the reference images, $recon_{ref}$. The middle columns show the object discovery results where each column corresponds to a single slot in $slots_{ref}$. The second rightmost column is the augmented images, img_{aug} . The rightmost column is the reconstruction of the augmented images, $recon_{aug}$. 11

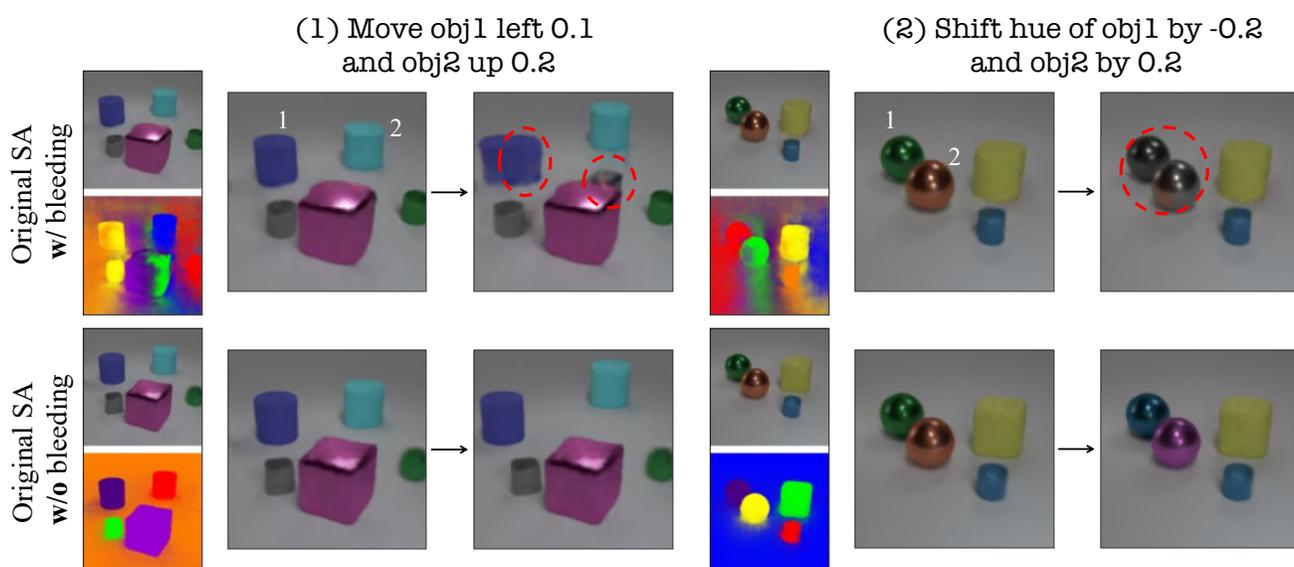


Figure 7. **Visualization of object manipulation results affected by the bleeding problem with the original Slot Attention.** The first row demonstrates the cases where bleeding problem emerges, while the second row shows the cases where the object discovery is done successfully.