

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

Aniket Didolkar^{1,2*} Andrii Zadaianchuk^{3*†} Rabiul Awal^{1,2*}
Maximilian Seitzer⁴ Efstratios Gavves^{3,5} Aishwarya Agrawal^{1,2†}
¹Mila - Quebec AI Institute, ²Université de Montréal
³University of Amsterdam, The Netherlands ⁴University of Tübingen
⁵Archimedes/Athena RC, Greece

Abstract

*Object-centric representation learning aims to decompose visual scenes into fixed-size vectors called “slots” or “object files”, where each slot captures a distinct object. Current state-of-the-art object-centric models have shown remarkable success in object discovery in diverse domains, including complex real-world scenes. However, these models suffer from a key limitation: they lack controllability. Specifically, current object-centric models learn representations based on their preconceived understanding of objects and parts, without allowing user input to guide which objects are represented. Introducing controllability into object-centric models could unlock a range of useful capabilities, such as the ability to extract instance-specific representations from a scene. In this work, we propose a novel approach for user-directed control over slot representations by conditioning slots on language descriptions. The proposed **CONTROLLABLE OBJECT-CENTRIC REPRESENTATION LEARNING** approach, which we term CTRL-O, achieves targeted object-language binding in complex real-world scenes without requiring mask supervision. Next, we apply these controllable slot representations on two downstream vision language tasks: text-to-image generation and visual question answering. We find that the proposed approach enables instance-specific text-to-image generation and also achieves strong performance on visual question answering.*

1. Introduction

Object-centric representation learning aims to decompose a visual scene into its constituent entities or objects and represent each entity as a distinct vector called a *slot*. Slot-based representations are inherently compositional and support many complex downstream tasks such as dynamics model-

ing [12, 43], control [6, 7, 16, 49], and reasoning [1, 29]. Moreover, studies in cognitive neuroscience [35, 41] have shown that human perception uses mechanisms akin to slot-based representations.

Existing unsupervised object-centric models [6, 9, 10, 28, 39] can successfully decompose complex real-world visual scenes. However, they face a fundamental limitation: they lack control over the object representations. While these models expose control over the *number* of scene parts, they do not allow users to extract specific object representations within a scene (e.g., specified by user queries in the form of language or position markers). For instance, a user can specify that a scene should be decomposed into K slots, but they cannot direct a given slot to bind to a particular object of interest such as “a cat” or “a black purse”.

This lack of control over the semantic content of the representation can be limiting, as this restricts the model always to extract a fixed decomposition of a scene based on its own preconceived understanding of objects and parts. Such rigidity can be problematic for applications that require representations at varying granularity, such as extracting the representation of a car wheel instead of the entire car, or vice versa.

Moreover, many downstream tasks may require or benefit from the knowledge of the semantic content in the slots. Due to the unsupervised nature of existing models, there is no way to identify the content of a slot without manually checking its corresponding mask. For example, if the user needs representations of the cat and a dog in a particular image to answer a question, they would have to manually inspect masks of all discovered objects to pick the corresponding slots.

To address these limitations, we propose to inject controllability into object-centric representation learning. We achieve this by querying the model to represent specific objects in the image. Specifically, these queries condition the slot vectors to guide them to the objects described by the query. The queries can be in the form of natural language (such as object category names or referring expressions).

*denotes equal contribution, order is determined by flipping coin

†denotes equal advising

The main challenge is to ensure that the slots conditioned on a specific query bind to the object referred to by that query. We term the challenge of binding slots to specific objects the *visual grounding problem* [15, 20, 34, 46]. We find that this problem is not trivial and introduce a novel controllable object-centric model — CTRL-O — to solve it. In our experiments, we demonstrate that the proposed approach can successfully bind slots to objects specified by user queries containing object categories or referring expressions in complex real-world scenes with limited supervision. In addition, we demonstrate the usefulness of the extracted controllable representations for two downstream tasks: visual question answering and instance-controllable image generation. Our contributions are as follows:

- We introduce CTRL-O, a novel method to learn controllable object-centric representations via user-defined inputs.
- We demonstrate that this approach supports extracting representations for complex reference expressions, enabling precise part specification within the representation.
- We validate the effectiveness of CTRL-O on two real-world downstream tasks: Instance-Controllable Image Generation and Visual Question Answering.

2. Related Works

Object-Centric Representation Learning Unsupervised object-centric representation learning (OCL) has gained a lot of interest in recent years [3, 5, 6, 9, 10, 14, 22, 28, 39, 50]. OCL aims to extract individual representations for various entities in unstructured sensory inputs such as images. Slot Attention [28] introduces an attention-based mechanism to decompose images into object-centric representations. DINO-SAU [39] builds upon this by utilizing self-supervised DINO features [4, 33] to enhance unsupervised object discovery. While DINOSAUR can effectively identify objects in real-world data [26], it lacks mechanisms for top-down control over the representations. In contrast, CTRL-O provides controllable OCL by incorporating language-based control queries, allowing for flexible guidance with minimal supervision during training. Some works [22, 23] have explored conditioning mechanisms in object-centric models. SAVi [22] uses bounding boxes for the initial frame of a video for conditioning. CoSA [23] conditions on learned vector representations. These methods are often limited to specific forms of conditioning and are primarily evaluated on synthetic datasets, while our method can handle many forms of conditioning on real-world data. Finally, several recent works connect object-centric representations with language [11, 21]. These works connect object representations with language post-hoc, assigning language labels to discovered slots. In contrast, CTRL-O integrates language and point conditioning directly into the learning process, allowing a user to control what representations should be extracted.

Downstream tasks with object-centric representations

There has been limited work exploring the applicability of object-centric models to downstream tasks. Slotformer [43] uses the learned slots for world modelling and video question answering. Zadaianchuk et al. [49], Yoon et al. [48] and Didolkar et al. [6] investigate the applicability of object-centric representations for learning RL policies in simple environments such as Atari [32]. One drawback of these works is that they mainly consider synthetic environments and toy tasks; thus, their applicability is limited. In contrast, in this paper, we consider downstream applications in complex real-world environments. There are only a few works that study applications of object-centric representations in real-world settings. Mamaghan et al [29] investigate the application of object-centric representations in Visual Question Answering. We consider them as a baseline for our experiments on VQA. Slot Diffusion [44] and Stable LSD [19] use object-centric representations for generating real-world images. However, both these approaches lack controllability; hence, it is difficult to specify any conditioning information or control the images that these approaches generate. In contrast, we demonstrate in Section 4.3, that CTRL-O, when used for image generation, provides fine-grained control over the image generation.

3. Method

In this section, we describe the proposed approach for injecting controllability into existing object-centric models. We present a visual depiction of our method in Fig. 1.

In our setup, the input consists of an image X and user-defined queries embedded into vectors $L = \{l_j \in \mathbb{R}^{D_{\text{emb}}}\}_{j=1}^M$. The expected object-centric representation of the image X is a set of slots $S = \{s_i \in \mathbb{R}^{D_{\text{slot}}}\}_{i=1}^N$. The first M slots (we assume that $M \leq N$) should represent the object identified by the corresponding queries, while the remaining slots represent the unspecified parts of the scene. This way, the obtained representation is a complete decomposition of the whole image X , while still containing the parts corresponding to the user-specified queries L .

We consider controllability in the form of *language queries*. We rely on the user to provide free-form text specifying object categories or object referring expressions whose visual representations are sought after. We encode this text into a fixed-sized vector embedding using LLM2Vec [2] with LLaMA-3-8B/LLaMA-3-8B [31] to obtain these embeddings. These language embeddings comprise the queries we feed into the model.

3.1. CTRL-O Architecture

Background We base the proposed approach on DINOSAUR [39]. DINOSAUR uses the Slot Attention module [28] for object discovery. Slot Attention is an attention-based differentiable clustering procedure which,

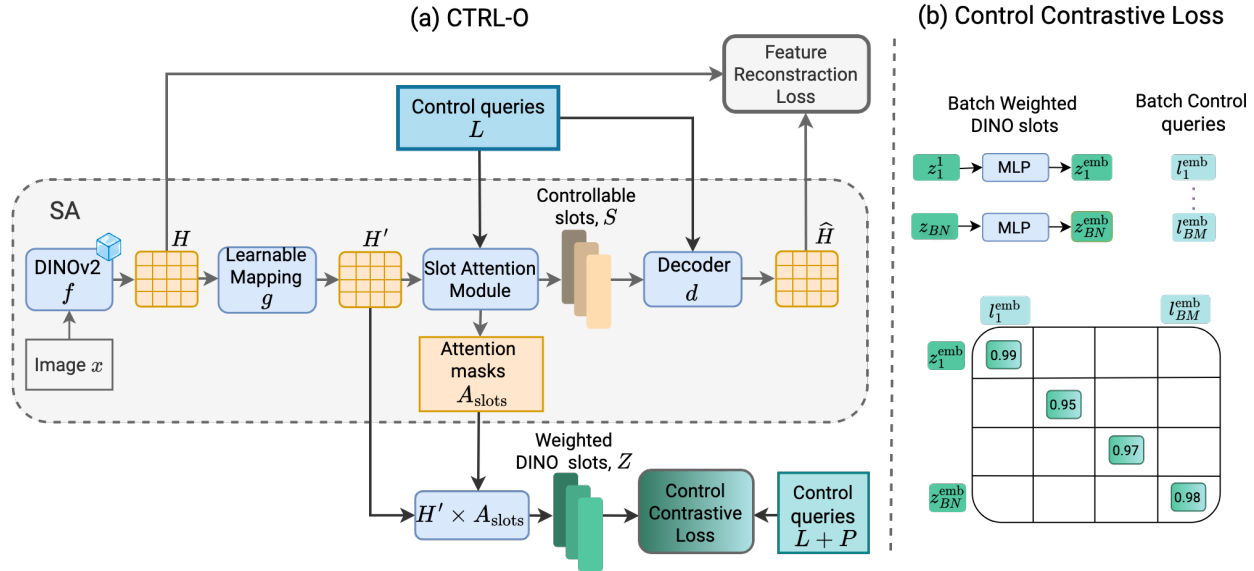


Figure 1. (a) Overview of CTRL-O architecture. An input image is processed by a frozen DINOv2 ViT model f , yielding patch features H . These features are then transformed into H' by a learnable transformer encoder g to align the feature space with the control queries. The control queries are introduced in the Slot Attention (SA) module, which guides the grouping of the encoded features into slots S . The initial slots in the SA module are conditioned with the control queries. Finally, an MLP decoder d , conditioned on control queries, reconstructs the DINOv2 features. (b) To ensure that slots utilize query information to represent specific objects, we apply a contrastive loss between control queries and the Slot Attention-modulated weighted DINO features A_{slot} (referred to as weighted DINO slots).

given a grid of features $H = \{h_k\}_{k=1}^K = f(X)$ obtained from an encoder f (we use DINOv2 [33]) applied to an image X , outputs a set of slots S such that each slot represents a distinct object in the image. We refer the reader to App. A for a detailed description of DINOSAUR.

Query-based Slot Initialization We are trying to solve the visual grounding problem: given the query corresponding to an object in the image, we want a slot to bind to exactly that object. The most straightforward way to enforce grounding is to condition the slots directly on the query corresponding to each object. Specifically, we achieve this by adding the object query l_i to one of the slots (see Fig. 1, input to the Slot Attention Module). This approach is similar to SAVi [22], which conditions each slot on the center of mass information for each object. In our experiments, we find that simply conditioning the slots on the queries does not lead to correct grounding; hence, a stronger signal is needed to ensure proper grounding.

Decoder Conditioning Similar to DINOSAUR, we use a broadcast MLP decoder, separately decoding each slot into patch features. We empirically find that conditioning the decoder on the corresponding control queries improves language grounding (concrete evaluation presented in Table 1. To implement this, we concatenate the resulting slots with

the control queries and pass them through an MLP whose output is fed into the patch decoder as shown in Fig. 1 (a).

3.2. Control Contrastive Loss to Enforce Grounding

To enforce grounding, we introduce a contrastive loss, as illustrated in Fig. 1 (b). The intuition behind this objective is that if a slot s_i is conditioned on a query l_i corresponding to the object o_i , then we want the encoder features corresponding to the slot s_i to be close in embedding space to the query l_i . To obtain the features corresponding to slot s_i , we spatially aggregate the features output by the mapping network (learnable mapping g in Fig. 1 (a)) by weighting them with the attention scores of slot s_i , obtained from the last iteration of slot attention: $z_i = \sum_{k=1}^K a_{ik} h_k$, where a_{ik} denotes the attention score of slot s_i on feature h_k . We further process z_i using an MLP to output z_i^{emb} , which is used in the contrastive loss.

Note that we do not directly use the slots for the contrastive loss because the loss can be trivially satisfied by the slots as they are conditioned on the control queries, which are the targets for the contrastive loss.

For the contrastive loss, we use (z_i^{emb}, l_i) as positive pairs which should be similar, and (z_i^{emb}, l_t) as negative pairs which should be dissimilar, with $t \neq i$. Note that for the negatives, we consider all the conditioning queries across the entire batch. Considering that there are T conditioning

| Slot Init. | GT Masks | CL | DC | Binding Hits | FG-ARI | mBO |
|------------|----------|----|----|--------------|-------------|-------------|
| ✓ | ✓ | ✗ | ✗ | 71.2 | 69.8 | 35.4 |
| ✓ | ✗ | ✗ | ✗ | 8.1 | 34.52 | 22.42 |
| ✓ | ✗ | ✗ | ✓ | 10.11 | 43.83 | 25.76 |
| ✓ | ✗ | ✓ | ✗ | 56.3 | 44.8 | 27.3 |
| ✓ | ✗ | ✓ | ✓ | 61.3 | 47.5 | 27.2 |

Table 1. **CTRL-O Model Component Ablation for Grounding.** Importance of various components for achieving strong grounding. Here, we use COCO *train* set for training and *val* set for evaluation. CL = Contrastive Loss, DC = Decoder Conditioning.

| Approach | Model | FG-ARI | mBO |
|-----------|------------------------------------|--------|------|
| Unsup. | DINOSAUR (MLP Dec.) [39] | 40.5 | 27.7 |
| | DINOSAUR (TF. Dec.) [39] | 34.1 | 31.6 |
| | Stable-LSD [19] | 35.0 | 30.4 |
| | SlotDiffusion [44] | 37.3 | 31.4 |
| Weak Sup. | Stable-LSD (Bbox Supervision) [40] | - | 30.3 |
| | CTRL-O (Trained on COCO) | 47.5 | 27.2 |

Table 2. **Object Discovery Performance.** Comparison of CTRL-O with unsupervised and weakly-supervised object-centric approaches on the COCO dataset.

queries in the entire batch, the loss for a single sample is formalized as follows:

$$\mathcal{L}_{CC}^l = - \sum_{i=1}^M \log \frac{\exp(z_i^{\text{emb}} \cdot l_i / \tau)}{\sum_{t=1}^T \exp(z_i^{\text{emb}} \cdot l_t / \tau)} \quad (1)$$

Here, τ is the temperature, which is set to 0.1. We assume two training regimes: when only language queries l_i are available and when both language queries l_i , and center-of-mass point queries p_i are provided during training. In the first regime, we define control contrastive loss \mathcal{L}_{CC} as simply \mathcal{L}_{CC}^l . In the second regime, control contrastive loss \mathcal{L}_{CC} is defined as the sum of two losses with language and point queries: $\mathcal{L}_{CC} = \mathcal{L}_{CC}^l + \mathcal{L}_{CC}^p$. We incorporate control contrastive loss \mathcal{L}_{CC} in addition to the feature reconstruction loss from DINOSAUR. For additional implementation details, see App. B.

4. Experiments

In this section, we first show that CTRL-O learns to bind to the right regions in the image given complex natural language queries. Next, we tackle two downstream tasks — instance-specific image generation and visual question answering — using a pretrained CTRL-O model.

4.1. Grounding Object-Centric Models

We study CTRL-O across two axes: 1) **Object Discovery** — How well can CTRL-O discover and represent each object separately in a scene?, and 2) **Grounding** — Can the slot conditioned on some language query l_j bind to the region specified by that language query?

Metrics To evaluate object discovery, we use standard metrics such as adjusted rand index (ARI) [18] and mean best overlap (mBO) [36]. To measure grounding, we introduce a new metric called *Binding Hits* which measures the grounding accuracy of the conditioned slots. Refer to App. F for more details regarding these metrics.

Datasets We use COCO [26] and Visual Genome [24] as our main datasets of study. COCO and Visual Genome both contain natural scenes with multiple objects. COCO contains category annotations spanning 91 different categories while Visual Genome contains region descriptions. COCO contains object annotations along with corresponding segmentation masks; we use it for quantitative evaluation of CTRL-O on object discovery and grounding. Visual Genome does not contain segmentation masks; hence, we only evaluate on it qualitatively. As several images in COCO contain multiple instances of the same object category, conditioning multiple slots on the same category name can be ambiguous for the model. Also, such conditioning poses problems for reliably computing the Binding Hits metric. Therefore, to disambiguate multiple instances of the same object category, we condition the slots on both category names and center of mass coordinates. We embed the language query using Meta-LLaMA-8B and the center of mass coordinates using a 2-layered MLP and concatenate them into a conditioning vector.

Object Discovery (Table 2) We compare CTRL-O to various unsupervised and one weakly-supervised object discovery method. All the methods considered in Table 2 apply Slot Attention to the features of a pretrained encoder to extract slots. Following DINOSAUR, this has become the standard in unsupervised object discovery. The weakly-supervised approach, Stable LSD (w/ bbox supervision) [40], uses bounding boxes to supervise the Slot Attention alpha masks. CTRL-O conditions the slots on language and center of mass queries and also uses the same information for contrastive loss. Therefore, CTRL-O can be classified as a weakly-supervised approach. Note that we do not use any of the guidance information to directly supervise the Slot Attention masks; thus, our form of supervision is weaker as compared to Stable LSD (w/ bbox supervision). In Table 2, we show that CTRL-O outperforms all unsupervised approaches in terms of ARI but lags behind in terms mBO. The lower performance in terms of mBO can be attributed to the MLP decoder of the underlying DINOSAUR model, which also obtains a lower mBO. We find that a transformer decoder [39] or a diffusion decoder [19, 40, 44] results in

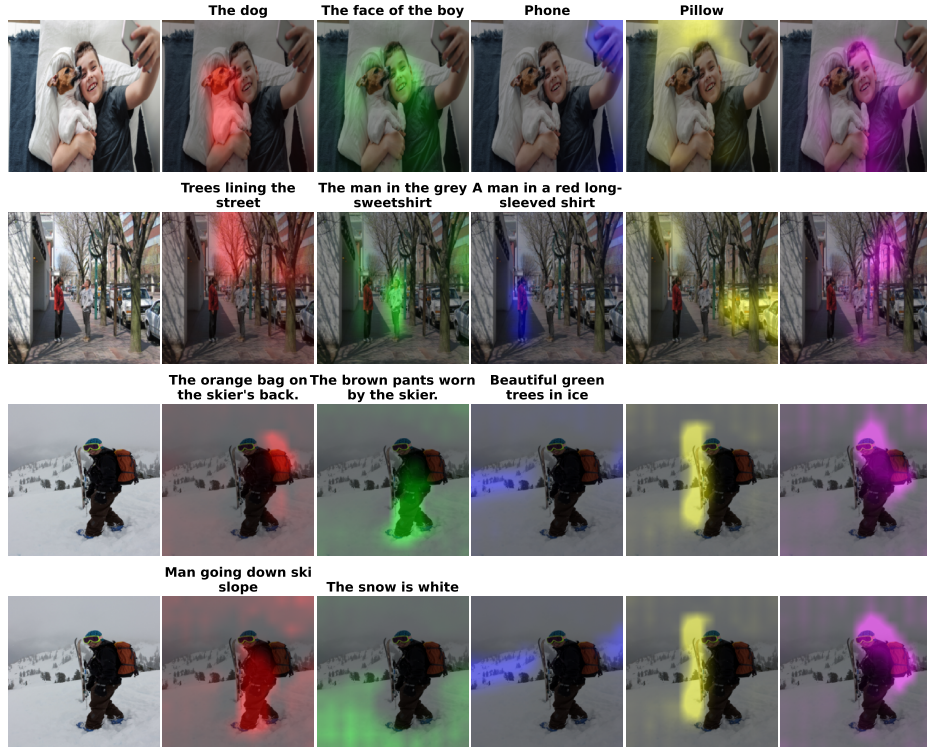


Figure 2. **Referring Expression Controllability on Visual Genome.** Visualization CTRL-O with free-form queries. The original image (left) and predicted segmentation masks are shown, with conditioning phrases presented above the corresponding segmented image; unconditioned slots have no phrase.

| Methods | Sup. | Image-text pretraining dataset | Fine tuning | RefCOCO | | | RefCOCO+ | | | Gref val |
|-----------------------|-----------------------------|--------------------------------|--------------|---------|-------|-------|----------|-------|-------|----------|
| | | | | val | testA | testB | val | testA | testB | |
| GroupViT [45] | \mathcal{T} | CC12M+YFCC | \times | 7.99 | 6.16 | 10.51 | 8.49 | 6.79 | 10.59 | 10.68 |
| | | | \checkmark | 10.82 | 11.11 | 11.29 | 11.14 | 10.78 | 11.84 | 12.77 |
| MaskCLIP [51] | \mathcal{T} | WIT | \times | 11.52 | 11.85 | 12.06 | 11.87 | 12.01 | 12.57 | 12.74 |
| | | | \checkmark | 19.45 | 18.69 | 21.37 | 19.97 | 18.93 | 21.48 | 21.11 |
| Shatter & Gather [21] | \mathcal{T} | VG | \times | 21.80 | 19.00 | 24.96 | 22.20 | 19.86 | 24.85 | 25.89 |
| CTRL-O | \mathcal{T} | VG | \times | 21.80 | 20.10 | 21.57 | 21.90 | 21.54 | 21.36 | 25.32 |
| CTRL-O | $\mathcal{T} + \mathcal{P}$ | VG | \times | 28.2 | 33.13 | 27.05 | 25.87 | 30.58 | 22.58 | 30.50 |

Table 3. **Referring expression segmentation** Comparison with weakly-supervised reference expression segmentation methods (Shatter & Gather) and open-vocabulary segmentation methods (GroupViT and MaskCLIP). The results on three datasets are reported in mIoU (%). Fine-tuning \checkmark means that the model is trained with the image-text pairs of the target benchmark; otherwise, the model is trained on the image-text pretraining dataset, and applied to the reference datasets zero-shot.

sharper masks as compared to the MLP decoder. This experiment verifies that CTRL-O can discover objects in complex natural scenes. Next, we evaluate whether it can bind to the region specified by the control queries.

Grounding Object Categories (Table 1) Controllability is a new paradigm for object-centric models that has not been explored before. Hence, there are no direct baselines with which we can compare. Instead, here we try to demonstrate

the difficulty of the grounding problem and ablate over the components introduced in Section 3 to understand their importance in achieving good grounding. We use the COCO dataset for this evaluation. Table 1 presents the results for various ablations. To obtain an upper bound for grounding performance, we train a (fully supervised) CTRL-O model by directly predicting the ground truth masks (first row in Table 1). While such supervised baseline achieves strong

segmentation performance (as indicated by ARI and mBO), it still cannot achieve perfect grounding (close to 100% Binding Hits), highlighting the difficulty of the grounding problem. Out of the components introduced in Sec. 3, the control contrastive loss is the most crucial component for achieving good grounding accuracy, followed by decoder conditioning. Without the contrastive loss, the model has no incentive to utilize the queries; hence, it does not achieve good Binding Hits values.

Grounding Referring Expressions (Figure 2) For COCO, we achieved controllability through both center of mass and category information. However, this approach is limited: COCO has a fixed number of categories, which affects generalizability. To overcome this issue, we can rely on *referring expressions*, where a user can refer to the target object using free-form natural language queries. To incorporate this ability, we use the Visual Genome dataset [24]. Since Visual Genome does not provide segmentation masks, we only evaluate CTRL-O qualitatively in this dataset.

We present the qualitative evaluation on Visual Genome in Fig. 2. We visualize the attention regions of each slot for a random sample of examples. We observe that slots conditioned on phrases bind to the object referred to in the phrase, while the unconditioned slots bind to the other objects not used for conditioning. The last two rows of the visualization demonstrate that CTRL-O can decompose the same scene at different levels of granularity, such as, just the leg (corresponding to the query “The brown pants worn by the skier”), just the bag (“The orange bag on the skier’s back”), or the entire skier (“Man going down ski slope”).

Referring Expression Segmentation Evaluation (Table 3) We evaluate CTRL-O trained on Visual Genome on referring expression segmentation on the RefCOCO, RefCOCO+, and Gref datasets. This is a zero-shot evaluation since these datasets were not used for training. We compare to various referring expression segmentation baselines, using mIoU between the predicted and ground truth mask as the metric. The most relevant baseline is Shatter & Gather (SaG) [21], which also employs slot attention to extract slots from the image, followed by cross attention to associate the referring expression with a slot. One limitation of SaG is that the referring expression does not directly influence the slot extraction process. This can be problematic for cases where the referring expression refers to an object not extracted by Slot Attention. This is not the case for CTRL-O, as in CTRL-O the referring expression directly influences the slot extraction. Moreover, all the baselines mentioned in Table 3 can only process a single referring expression per forward pass, while CTRL-O can process multiple referring expressions in parallel by conditioning multiple slots on the corresponding expressions. We find that CTRL-O outperforms all the baselines. However, CTRL-O by default uses language queries and center of mass annotations ($\mathcal{T} + \mathcal{P}$)

while the baselines use only language queries as weak supervision. Hence, we also run CTRL-O with only language queries. Implementation details for this variant are provided in App. C. We find that CTRL-O with only language queries (\mathcal{T}) achieves competitive performance with SaG.

4.2. Unlocking New Capabilities for Object-Centric Models with CTRL-O

In this section, we show that CTRL-O can be used for *instance controllable image generation* (Fig. 3(a)) — a use case where existing object-centric methods fail. We also demonstrate how CTRL-O can be used in a novel way to improve over existing object-centric methods in Visual Question Answering (VQA) (Fig. 3(b)). In these experiments, we use CTRL-O pre-trained on both Visual Genome and COCO with language-only conditioning (referring expressions in Visual Genome and Object Categories in COCO).

4.3. Instance Controllable Image Generation

In this section, our goal is to demonstrate that object-centric representations obtained from CTRL-O can be used for controllable image generation. Specifically, we aim for control in the space of instances where specific instances can be extracted from images and used as a conditioning for generating images containing those instances. Prior works, such as SlotDiffusion [44] and Stable-LSD [19], use Stable Diffusion [38] as a decoder to reconstruct images from slot vectors. These models condition the U-Net in Stable Diffusion on slots obtained from Slot Attention, enabling slot-conditional generation. However, these approaches are fundamentally limited: 1) To control the instances in the generated image, the user needs to manually reconstruct the masks corresponding to all the slots and find target slots that correspond to instances of interest. 2) They fix the diffusion model to a fixed number of slots, as in Slot Attention, limiting flexibility. Users cannot condition on a subset of objects while leaving the image layout flexible, and text inputs are unsupported.

CTRL-O addresses the above limitations: 1) Language-based control is inherent to CTRL-O. Hence, to control the instances present in the generated image, a user needs only to query CTRL-O to extract the corresponding instance representations from a given image as shown in Fig. 3(a). 2) We maintain the text interface of Stable Diffusion and only add the slots to control the visual identity of specific instances. This is similar to [47] where a user can specify an image containing an object for instance controlled generation. Here instead of specifying images, we use slots obtained from CTRL-O for this task. For example, prompting Stable Diffusion with “A photo of a bus” generates a random bus. But if we want a specific bus as in Fig. 3(a), we use CTRL-O to extract its representation from an image, allowing us to prompt the diffusion model with “A photo of a bus. S_{bus} ”,

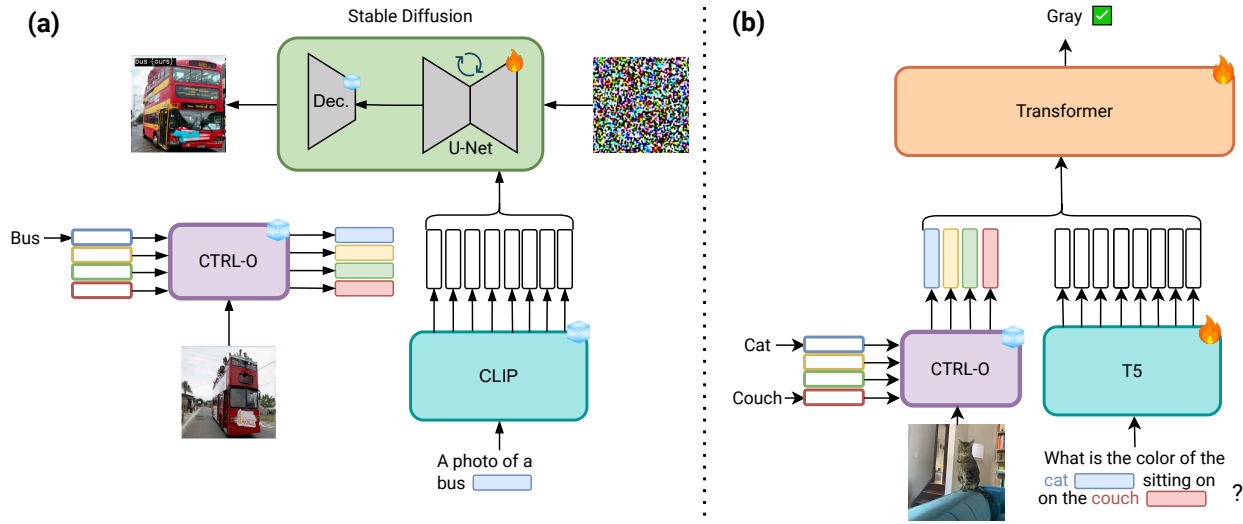


Figure 3. (a) **Instance Specific Image Generation**: we query an image to extract slot representations of instances required in the generated image. The slot corresponding to that instance is then fed into a Stable Diffusion model along with the caption to generate the corresponding image. (b) **Visual Question Answering**: we embed the slots directly into the question. Given a question, we parse it into a set of noun chunks or referring expressions and use them to extract the corresponding representations with the CTRL-O. The slots are then embedded into the appropriate positions into the text and fed into the language model. 🧊 = frozen parameters; 🔥 = trainable parameters.

where S_{bus} is obtained from CTRL-O.

CTRL-O-SD A visual depiction of our generation pipeline is presented in Fig. 3(a). We keep CTRL-O frozen. We use stable diffusion as the generative model [38]. We finetune the U-Net in Stable Diffusion while keeping the rest of the components fixed. The COCO dataset is used for finetuning. We feed the model captions as conditioning inputs. We take the category names present in the caption and extract the corresponding slots using CTRL-O. We then project the slots into CLIP embeddings space and append them to the language caption as shown Fig. 3(a). We refer to the proposed pipeline as CTRL-O-Stable Diffusion (CTRL-O-SD). Note that this is different to the approach adopted in Stable LSD where the slot attention module is trained in conjunction with the diffusion model. Therefore, we treat generation as a downstream task for CTRL-O while Stable LSD uses diffusion to train the object-centric model.

Results (Figure 4, Table 4, and Table 5) We first compare CTRL-O-SD to an existing object-centric method - Stable LSD - in terms of generative quality as measured by Fréchet Inception Distance (FID). Note that we can only compare CTRL-O-SD to Stable LSD in terms of image generation quality since Stable LSD lacks controllability and hence cannot be compared to CTRL-O-SD along that axis. We find that CTRL-O-SD achieves higher generation quality than Stable LSD as evidenced by the lower FID in Table 4. Next, we evaluate CTRL-O-SD along the axis of instance controllability in the generated output on the COCO

validation set. Here we use Stable Diffusion for comparison. For Stable Diffusion, we feed the captions as conditioning while for CTRL-O-SD, we feed the captions along with slots extracted from the ground truth images from the COCO validation set by conditioning them on the categories present in the caption. We report the CLIP-I Score metric, which measures the cosine similarity between the CLIP vision embedding of the generated and the ground truth images (see further details in App. H.2). Table 5 shows that CTRL-O-SD achieves a higher CLIP-I Score, thus exhibiting stronger instance controllability.

| Method | FID Score (↓) |
|----------|---------------|
| LSD [19] | 26.20 |
| CTRL-O | 25.20 |

Table 4. **Stable LSD Comparison** We compare CTRL-O and Stable LSD in terms of image generation quality. We find that CTRL-O achieves a lower FID as compared to LSD.

We also visualize the generated images from CTRL-O-SD and SD in Figure 4. We can see that the images generated using CTRL-O-SD contain instances that are closer to the query image. We also show that CTRL-O-SD can compose slots corresponding to two different categories “A bench” and “A pizza” from two different images, to produce an image specified by “A pizza on a bench. $S_{Bench} S_{Pizza}$ ”. We showcase some failure modes of CTRL-O-SD in App. H.4.

a. Instance Controllable Image Generation



b. Multi-Instance Composition



Figure 4. **a. Instance Controllable Image Generation.** Comparison between CTRL-O-SD and the baseline Stable Diffusion (SD). For a given query image (marked **query**), we extract a slot representation of a specific instance I_q (e.g., laptop, bus, banana). In CTRL-O-SD, the input is “A photo of I_q . S_{I_q} ” to guide instance generation, while for SD, only “A photo of I_q ” is used. Our approach produces images that more closely match the visual identity of the conditioned instance **b. Multi-Instance Composition.** We extract instances from multiple images (e.g., “bench” and “pizza”) and compose them into a single image, as seen with “the pizza on the bench”.

In these cases, CTRL-O struggles to bind to the correct object specified by the query, and the diffusion model exhibits issues with object deformation and repetitions.

4.4. Visual Questions Answering

Here we consider the task of Visual Question Answering on the VQA2 dataset [13]. We treat this problem as a classification problem and use accuracy as the metric. Language-guided object-centric representations from CTRL-O can potentially provide a much stronger coupling between the vision and language inputs in VQA. To achieve this, we

propose to directly insert the slots into the question before feeding it into the language model. The proposed approach is presented in Fig. 3(b). Given a question, we first use spaCy to extract noun chunks (e.g., Cat, Couch) for the image, which are then used to condition the slots in CTRL-O to extract the corresponding slots. These slots are then inserted into the question at the appropriate positions as shown in Fig. 3(b). For example, the question “What is the color of the cat sitting on the couch?” becomes “What is the color of the cat S_{cat} sitting on the couch S_{couch} ?”. We use a learnable linear projection to map the slots to the same di-

| Method | CLIP-I Score (\uparrow) |
|---------|-----------------------------|
| SD [38] | 0.71 |
| CTRL-O | 0.78 |

Table 5. **Stable Diffusion Comparison** We compare CTRL-O to Stable Diffusion on Instance Controllable Image Generation. CLIP-I score measures the cosine similarity in the CLIP Image embedding space between the generated and the query image. We can see that CTRL-O achieves a higher CLIP-I score, showing that the generated instances are closer to those in the query image.

mension as the T5 embeddings. We feed this question, with the inserted slots, into the language model. Therefore, the language and vision inputs are strongly coupled from the input stage, which allows for more interaction between visual and language components. We refer to this technique as *coupling*. This technique is reminiscent of adding additional learnable tokens as input to the language model [25].

| Model | No-Coupling | Coupling |
|----------------|-------------|----------|
| DINOv2 (22M) | 58.26 | 58.12 |
| CLIP (191M) | 58.43 | 58.64 |
| DINOSAUR (37M) | 58.32 | 57.66 |
| CTRL-O (39M) | 59.18 | 60.25 |

Table 6. **VQA performance on VQAv2.** We report the standard VQA accuracy metric along with the number of parameters per model (in parentheses). Our proposed (CTRL-O) achieves the highest accuracy in both No-Coupling and Coupling settings, demonstrating its effectiveness.

Setup Our full pipeline is shown in Fig. 3(b), inspired by [29]. We use T5 as the language embedding model, and the vision model is CTRL-O. For the baselines, we use CLIP [37], DINOv2 [33], and DINOSAUR [39] as the vision models. The output network is a transformer with 2 layers and 64 heads. For all methods, we feed visual representations (patches for CLIP and DINOv2 and slots for CTRL-O and DINOSAUR) and the language embeddings from T5 into the output network. The output network is trained from scratch, while T5 undergoes fine-tuning.

We introduce two variants of VQA training with CTRL-O: 1) CTRL-O that directly feeds the slots into the Transformer output network without embedding them in the language, similar to baseline methods; 2) CTRL-O (with coupling) that inserts the corresponding slot representations into the appropriate place in the question as shown in Fig. 3(b). To ensure a fair comparison, we extend coupling experiments to all baselines. Since these models lack explicit object binding, we insert their aggregated features (CLS token for DINOv2/CLIP and slot mean for DINOSAUR) into the same positions as CTRL-O’s control slots. This allows us to assess how different representations interact

with language. Classification accuracy is reported (see details in Appendix G).

Results Table 6 shows that both CTRL-O variants outperform all baselines, demonstrating the effectiveness of its object-centric representations. Notably, coupling primarily benefits CTRL-O, as its representations explicitly align with language. In contrast, baselines insert generic image representations that do not naturally correspond to the preceding text, limiting the effectiveness of coupling. However, while our method improves performance, it remains below the latest state-of-the-art models (>80%) [27], which leverage large language models (LLMs) and web-scale multimodal data.

5. Conclusion

We have introduced CTRL-O, a controllable object-centric model that can be queried to extract representations of specific objects in a scene. We experimentally showed that representations of specific objects can be extracted in complex real-world scenes based on a range of user queries such as object category names and referring expressions. This capability expands the applicability of object-centric models to various real-world applications, such as instance-controllable image generation and visual question answering. Therefore, our work takes a step towards improving the applicability of object-centric models to complex real-world downstream tasks. Future work should explore how to improve the grounding of these models. The current approach, while useful, is still limited since it does not achieve perfect grounding. Moreover, we hope that learning controllable object-centric representations becomes the standard way of learning object-centric representations and leads to broader adoption of object-centric models to various domains and downstream tasks.

Acknowledgments

The authors would like to thank Kanishk Jain, Ankur Sikarwar, and Le Zhang for reviewing and providing feedback on an earlier version of the paper and Amir M. Karimi and Samuele Papa for providing the code for the VQA experiments in Section 4.4. This research was enabled in part by compute resources provided by Mila (mila.quebec). AD would like to thank Nanda Harishankar Krishna, Moksh Jain, and Rohan Banerjee for help with Fig. 3. Additionally, AD would like to thank Anirudh Goyal and Mike Mozer for helpful discussions regarding the research direction explored in this paper. AD also acknowledges the support of a scholarship from UNIQUE (<https://www.unique.quebec>). AZ is funded by the European Union (ERC, EVA, 950086)

Contributions

AD and RA initially started a collaboration on object-centric models for vision-language tasks. AD came up with the idea of controllable slots conditioned on language. AZ proposed the contrastive objective for the model. AD and AZ developed the code for the model and conducted experiments on object grounding. AD ran all the baselines in Section 4.1. RA developed the code and ran experiments for instance-controllable image generation. AZ and RA developed the initial code for the VQA experiment, which AD and RA further extended for the final results. AZ, RA, and AD ran the baselines for Section 4.4. AA and AZ provided supervision and guidance throughout the project. MS and EG took part in discussions. AD, RA, AZ, and AA wrote the paper with help from MS.

References

- [1] Rim Assouel, Pau Rodriguez, Perouz Taslakian, David Vazquez, and Yoshua Bengio. Object-centric compositional imagination for visual abstract reasoning. In *ICLR2022 Workshop on the Elements of Reasoning: Objects, Structure and Causality*, 2022. 1
- [2] Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. Llm2vec: Large language models are secretly powerful text encoders. *arXiv preprint arXiv:2404.05961*, 2024. 2, 1
- [3] 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. 2
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging Properties in Self-Supervised Vision Transformers. *ICCV*, 2021. 2, 1
- [5] Aniket Didolkar, Andrii Zadaianchuk, Anirudh Goyal, Mike Mozer, Yoshua Bengio, Georg Martius, and Maximilian Seitzer. Zero-shot object-centric representation learning. *arXiv preprint arXiv:2408.09162*, 2024. 2
- [6] Aniket Rajiv Didolkar, Anirudh Goyal, and Yoshua Bengio. Cycle consistency driven object discovery. In *ICLR*, 2024. 1, 2
- [7] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. PaLM-E: An embodied multimodal language model. In *ICML*, 2023. 1
- [8] Constantin Eichenberg, Sidney Black, Samuel Weinbach, Letitia Parcalabescu, and Anette Frank. Magma—multimodal augmentation of generative models through adapter-based finetuning. *arXiv preprint arXiv:2112.05253*, 2021. 4
- [9] Martin Engelcke, Adam R. Kosiorok, Oiwi Parker Jones, and Ingmar Posner. GENESIS: Generative Scene Inference and Sampling with Object-Centric Latent Representations. In *ICLR*, 2020. 1, 2
- [10] S. M. Ali Eslami, Nicolas Heess, Theophane Weber, Yuval Tassa, David Szepesvari, Koray Kavukcuoglu, and Geoffrey E. Hinton. Attend, Infer, Repeat: Fast Scene Understanding with Generative Models. In *NeurIPS*, 2016. 1, 2
- [11] Ke Fan, Zechen Bai, Tianjun Xiao, Dominik Zietlow, Max Horn, Zixu Zhao, Carl-Johann Simon-Gabriel, Mike Zheng Shou, Francesco Locatello, Bernt Schiele, Thomas Brox, Zheng Zhang, Yanwei Fu, and Tong He. Unsupervised open-vocabulary object localization in videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13747–13755, 2023. 2, 1
- [12] Anirudh Goyal, Aniket Didolkar, Nan Rosemary Ke, Charles Blundell, Philippe Beaudoin, Nicolas Heess, Michael C Mozer, and Yoshua Bengio. Neural production systems. *Advances in Neural Information Processing Systems*, 34:25673–25687, 2021. 1
- [13] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017. 8
- [14] 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. In *ICML*, 2019. 2
- [15] Klaus Greff, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. On the Binding Problem in Artificial Neural Networks. *arXiv:2012.05208*, 2020. 2
- [16] Dan Haramati, Tal Daniel, and Aviv Tamar. Entity-centric reinforcement learning for object manipulation from pixels. In *ICLR*, 2024. 1
- [17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a nash equilibrium. *CoRR*, abs/1706.08500, 2017. 5
- [18] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*, 1985. 4, 2
- [19] Jindong Jiang, Fei Deng, Gautam Singh, and Sungjin Ahn. Object-centric slot diffusion. In *NeurIPS*, 2023. 2, 4, 6, 7
- [20] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. ReferItGame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 787–798, Doha, Qatar, 2014. Association for Computational Linguistics. 2
- [21] Dongwon Kim, Namyup Kim, Cuiling Lan, and Suha Kwak. Shatter and gather: Learning referring image segmentation with text supervision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15547–15557, 2023. 2, 5, 6, 1
- [22] Thomas Kipf, Gamaleldin Fathy Elsayed, Aravindh Mahendran, Austin Stone, Sara Sabour, Georg Heigold, Rico Jon-

- schkowsky, Alexey Dosovitskiy, and Klaus Greff. Conditional Object-centric Learning from Video. In *ICLR*, 2022. 2, 3
- [23] Avinash Kori, Francesco Locatello, Francesca Toni, and Ben Glocker. Unsupervised conditional slot attention for object centric learning. *arXiv preprint arXiv:2307.09437*, 2023. 2
- [24] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. Visual genome: Connecting language and vision using crowdsourced dense image annotations, 2016. 4, 6
- [25] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, 2021. 9
- [26] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common Objects in Context. In *ECCV*, 2014. 2, 4
- [27] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 9
- [28] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-Centric Learning with Slot Attention. In *NeurIPS*, 2020. 1, 2
- [29] Amir Mohammad Karimi Mamaghan, Samuele Papa, Karl Henrik Johansson, Stefan Bauer, and Andrea Dittadi. Exploring the effectiveness of object-centric representations in visual question answering: Comparative insights with foundation models. *arXiv:2407.15589*, 2024. 1, 2, 9, 4
- [30] Oscar Mañas, Pau Rodriguez, Saba Ahmadi, Aida Nematzadeh, Yash Goyal, and Aishwarya Agrawal. Mapl: Parameter-efficient adaptation of unimodal pre-trained models for vision-language few-shot prompting. *arXiv preprint arXiv:2210.07179*, 2022. 4
- [31] Meta. The llama 3 herd of models, 2024. 2
- [32] Volodymyr Mnih. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013. 2
- [33] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafranec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *TMLR*, 2023. 2, 3, 9, 1
- [34] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, 2018. 2
- [35] Steven Pinker. Visual cognition: An introduction. *Cognition*, 1984. 1
- [36] Jordi Pont-Tuset, Pablo Arbeláez, Jonathan T. Barron, Ferran Marques, and Jitendra Malik. Multiscale Combinatorial Grouping for Image Segmentation and Object Proposal Generation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. 4, 2
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 9, 1
- [38] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 6, 7, 9
- [39] Maximilian Seitzer, Max Horn, Andrii Zadaianchuk, Dominik Zietlow, Tianjun Xiao, Carl-Johann Simon-Gabriel, Tong He, Zheng Zhang, Bernhard Schölkopf, Thomas Brox, and Francesco Locatello. Bridging the gap to real-world object-centric learning. In *ICLR*, 2023. 1, 2, 4, 9
- [40] Krishnakant Singh, Simone Schaub-Meyer, and Stefan Roth. Guided latent slot diffusion for object-centric learning, 2024. 4
- [41] Elizabeth S. Spelke. Core knowledge. *The American psychologist*, 2000. 1
- [42] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 4
- [43] Ziyi Wu, Nikita Dvornik, Klaus Greff, Thomas Kipf, and Animesh Garg. SlotFormer: Unsupervised visual dynamics simulation with object-centric models. In *ICLR*, 2023. 1, 2
- [44] Ziyi Wu, Jingyu Hu, Wuyue Lu, Igor Gilitschenski, and Animesh Garg. Slotdiffusion: Object-centric generative modeling with diffusion models. In *NeurIPS*, 2023. 2, 4, 6
- [45] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18134–18144, 2022. 5
- [46] Kelvin Xu. Show, attend and tell: Neural image caption generation with visual attention. *arXiv preprint arXiv:1502.03044*, 2015. 2
- [47] Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. Paint by example: Exemplar-based image editing with diffusion models. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18381–18391, 2022. 6
- [48] Jaesik Yoon, Yi-Fu Wu, Heechul Bae, and Sungjin Ahn. An investigation into pre-training object-centric representations for reinforcement learning. In *ICML*, 2023. 2
- [49] Andrii Zadaianchuk, Maximilian Seitzer, and Georg Martius. Self-supervised Visual Reinforcement Learning with Object-centric Representations. In *ICLR*, 2020. 1, 2
- [50] Andrii Zadaianchuk, Maximilian Seitzer, and Georg Martius. Object-centric learning for real-world videos by predicting temporal feature similarities. In *NeurIPS*, 2023. 2
- [51] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *European Conference on Computer Vision (ECCV)*, 2022. 5

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

Supplementary Material

A. DINOSAUR Implementation Details

DINOSAUR uses a DINO [4] encoder to process the image into features. It relies on a feature reconstruction loss to supervise the object discovery process. Throughout the training, the DINO encoder is kept frozen. We adopt a similar approach, however we use a DINOv2 [33] encoder instead of a DINO encoder. Figure 5 illustrates the DINOSAUR architecture with a DINOv2 backbone. Additionally, we have added a learnable mapping network g , which is a 3-layer Transformer after the frozen DINOv2 encoder. SA module is applied on top of the mapping output as shown in Figure 1(a).

B. CTRL-O Implementation Details

Algorithm 1 Slot Attention with Language Conditioning

Input: $\text{inputs} \in \mathbb{R}^{N \times D_{\text{inputs}}}$, $\text{slots} \in \mathbb{R}^{K \times D_{\text{slot}}}$, language queries $\ell \in \mathbb{R}^{M \times D_{\text{lang}}}$

Layer params: k, q, v : linear projections for attention; p_ℓ : projection for language query; GRU; MLP; LayerNorm (x3)

```
1:  $\text{inputs} \leftarrow \text{LayerNorm}(\text{inputs})$ 
2:  $\ell_{\text{proj}} \leftarrow p_\ell(\ell)$   $\triangleright$  Project  $M$  language queries to slot dimension
3:  $\{\text{slots}\}_{i=1}^M \leftarrow \ell_{\text{proj}}$   $\triangleright$  Condition first  $M$  slots on language query
4: for  $t = 0 \dots T - 1$  do
5:    $\text{slots}_{\text{prev}} \leftarrow \text{slots}$ 
6:    $\text{slots} \leftarrow \text{LayerNorm}(\text{slots})$ 
7:    $\text{attn} \leftarrow \text{Softmax}(\frac{1}{\sqrt{D}} k(\text{inputs}) \cdot q(\text{slots})^\top \text{axis} = \text{slots})$ 
8:    $\text{updates} \leftarrow \text{WeightedMean}(\text{weights} = \text{attn} + \epsilon, \text{values} = v(\text{inputs}))$ 
9:    $\text{slots} \leftarrow \text{GRU}(\text{state} = \text{slots}_{\text{prev}}, \text{inputs} = \text{updates})$ 
10:   $\text{slots} \leftarrow \text{slots} + \text{MLP}(\text{LayerNorm}(\text{slots}))$ 
11: end for
12: return  $\text{slots}$ 
```

We present the modified Slot Attention with query-based initialization in Algorithm 1.

Control Contrastive Loss For conditioning, we mainly use language queries. However, we assume that each image in our dataset consists of multiple object annotations, each containing a center of mass annotation and a category or

referring expression annotation. Therefore, we have two separate contrastive losses - one each for the language information and the point information, as shown in Figure 1(b).

Conditioning We run Slot Attention for a fixed number of slots K . However, in general, we may not have K queries per image. In such cases, we initialize a subset of the slots with the given queries, and the rest are free to bind to any of the other objects in the scene (see line 3 of Algorithm 1). When computing the contrastive loss, we only consider slots conditioned on some query.

C. Training CTRL-O with Language Queries

Needing center of mass annotations for the contrastive loss can be a limitation as these annotations may not be available in many datasets. Further, the main baseline that we consider for the referring expression segmentation task (Section 4.1) - Shatter-and-Gather [21] - does not require center of mass annotations. Therefore, for an apples-to-apples comparison, we implement a variant of CTRL-O which does not require center of mass annotations.

A visual depiction of this approach is presented in Figure 7. First, we remove additional center-of-mass information and leave only language queries in the contrastive loss. We find that simply removing the center-of-mass information leads to collapse of representations as the contrastive loss can be trivially satisfied by directly using the language embeddings on which the slots are conditioned on - we term this as *leakage*. To prevent leakage we propose to use CLIP [37] image features and language embeddings in the control contrastive loss. In particular, instead of taking the weighted average of DINO features (Figure 1(a)), we take the weighted average of patch-based CLIP features [11]. The slot conditioning still uses language embeddings from LLaMa-3-8B [2], however, CLIP language embeddings are used as targets in the contrastive loss. This way, CTRL-O learns to bind to the correct regions in the image specified by language queries without center-of-mass annotations.

D. Choice of Decoder for CTRL-O

In this subsection, we additionally study the compatibility of CTRL-O with different previously proposed decoders. In particular, we investigate the compatibility and scalability of our method with two different decoder architectures (MLP and Transformer). In Table 2, we compare our method with other OCL methods, showing that while our method strongly

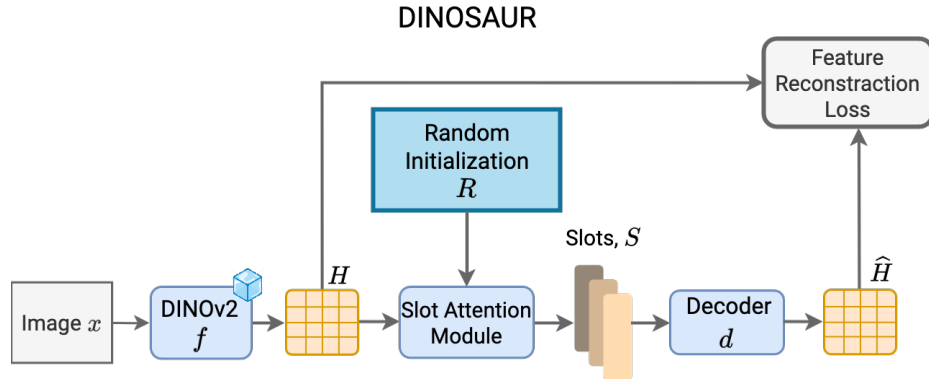


Figure 5. Overview of DINOSAUR architecture. The image is processed into a set of patch features H by a frozen DINO ViT model. The Slot Attention module groups the encoded features into a set of slots initialized by random queries sampled from the same Gaussian distribution with learnable parameters. By contrast, CTRL-O is initialized by the combination of control queries for conditioned slots and random queries for unconditioned slots. DINOSAUR is trained by reconstructing the DINO features from the slots using MLP decoder [39].

outperforms other methods in FG-ARI, its mask quality is lower than methods with stronger (pretrained) diffusion and Transformer decoders that have less inductive bias towards scene decomposition. Thus, it is important to investigate how our method performs with different decoders and whether we can scale MLP decoders for better mask quality. Object discovery with the Transformer decoder was shown to be sensitive to hyperparameters and can entirely fail (see App. D.4 and D.5 of DINOSAUR paper [39]). Subsequently, we also find that CTRL-O with Transformer decoder achieves 10.2 mBO. Through thorough investigation, we conclude that *Transformer decoder is not compatible with contrastive loss, which is needed for language controllability in CTRL-O but not in the baselines*. Thus, to improve masks quality we propose to scale the MLP decoder itself; scaling MLP dim to 4096 led to improved 28.0 mBO and 47.9 FG-ARI.

E. Referring Expression Visualization

In Figure 6, we compare the visualizations obtained from CTRL-O ($\mathcal{L}+\mathcal{P}$ setting), CTRL-O (\mathcal{L} setting), and Shatter-and-Gather (SaG). Note that the queries listed on the top of each column are free-form queries created by a user and may not be similar to those typically found in the visual genome dataset. One potential issue with Shatter-and-Gather is that the language queries do not influence the slot extraction process - Slot attention first extracts a fixed number of slots, after which the query binds to the most relevant slot post-hoc. This can be limiting, as in some cases, the region referred to by the query may not be extracted into a single slot. In such cases, the language query may not bind to any slot. In Figure 6, we find that this is exactly what happens in several cases for Shatter-and-Gather. For example, in the first column, for the query “The orange bag on the skier’s bag”, SaG binds to the skier’s shoes. In the 5th column, SaG fails

to bind to any region for the query “the lamp”. In contrast, both variants of CTRL-O frequently bind to the correct regions specified by the queries. Secondly, in CTRL-O the language queries directly influence slot extraction which allows it to explicitly extract the referred regions from the image and bind to them.

A particularly interesting case is the last row for CTRL-O (\mathcal{L}), where it learns to bind correctly even though queries are less specific and more subjective - “the ancient building” and “the new building”. This emphasizes the generalizability of CTRL-O to complex language queries.

F. Object Discovery and Binding Metrics

FG-ARI The *adjusted rand index* (ARI) measures the similarity between two clusterings [18]. We use the instance/object masks as the targets. We only compute this metric for pixels in the foreground (hence, FG-ARI). Unlabeled pixels are treated as background.

mBO To compute the mBO [36], each ground truth mask (excluding the background mask) is assigned to the predicted mask with the largest overlap in terms of IoU. The mBO is computed as the average IoU of these mask pairs.

Binding Hits This metric measures controllable grounding. For binding hits, consider that a slot s_i is conditioned on a query L_i identifying an object o_i with ground-truth mask m_i . The broadcast decoder of slot attention outputs a mask per slot. If the overlap between the predicted mask for slot s_i , denoted as \hat{m}_i , and the ground truth mask m_i is the highest among all pairs of predicted and ground truth masks, it is considered as a hit (1) else it is considered as a miss (0). Binding Hits metric is measured as the average number of hits across the dataset.

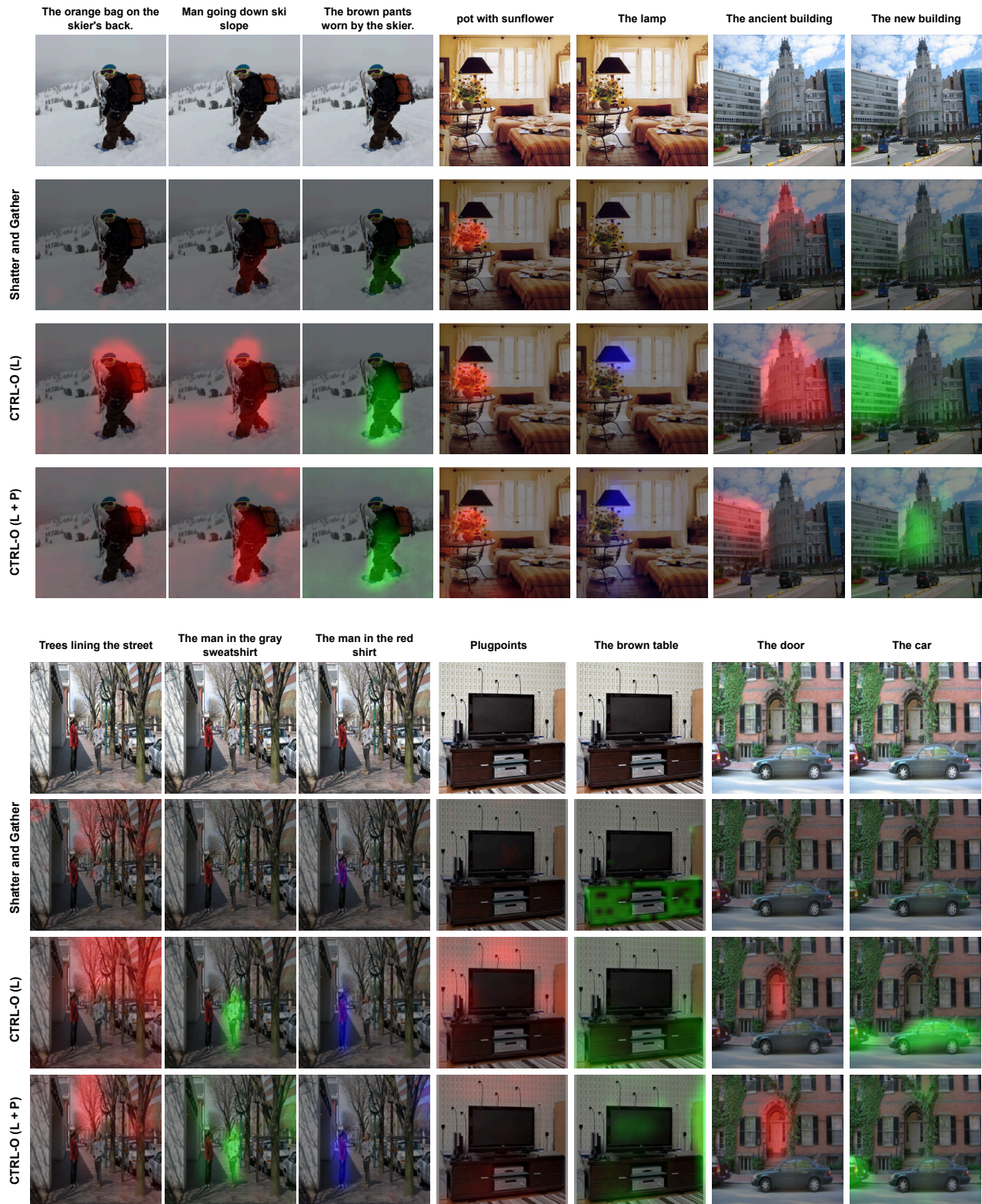


Figure 6. **Visualization Comparison** In this figure we visualize and compare the masks obtained using CTRL-O ($\mathcal{L} + \mathcal{P}$), CTRL-O (\mathcal{L}), and SaG when queried with free-form language queries. We can see that both the variants based on CTRL-O are significantly better at binding to the correct region descriptions as compared to SaG. This difference can be attributed CTRL-O using the language guidance to directly influence the slot extraction process while SaG considers the language to slot binding as a post-processing step after the slots have been extracted.

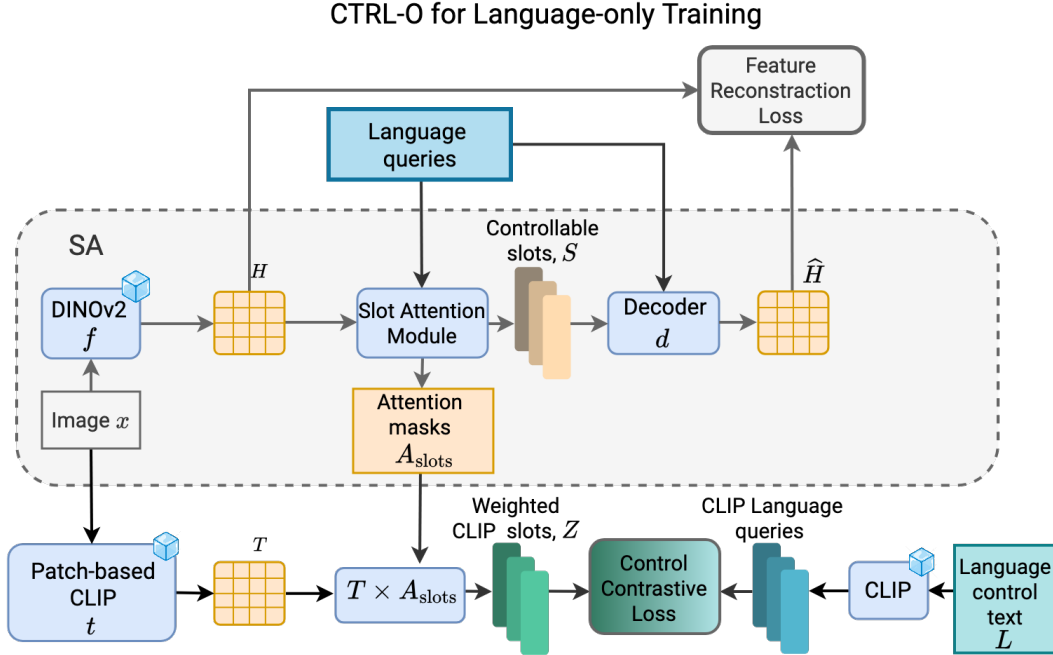


Figure 7. **Language-Only CTRL-O** Overview of language-only training. In this setting, we use the frozen CLIP model to compute both weighted CLIP slots and CLIP Language queries that we use in control contrastive loss. We average features from CLIP using attention weights from the Slot Attention module.

G. Additional Details of VQA Experiments

Evaluation Metric. We evaluate VQA models using classification accuracy across 3000 classes, using the top-frequent answers, which covers more than 90% of the question in the dataset.

Discussion on Coupling in CTRL-O VQA model The standard approach for solving VQA tasks with pretrained vision and language backbones is to feed the output representations of the vision model and the language model into a single neural network - usually a Transformer [42] - which then outputs a distribution over the answer categories [8, 29, 30]. To solve VQA, it is crucial to have strong interaction between the visual and language inputs. However, in pre-existing approaches, this interaction only happens in the output network (the Transformer that processes the language and vision outputs), which can be limiting.

To address this, we introduce an approach called *coupling*. Coupling, with the help of CTRL-O, directly inserts the visual representations into the language query, thus enabling strong vision and language interaction from the input stage. The proposed approach is presented in Figure 3(b).

H. CTRL-O SD

H.1. Fine-Tuning Details

In CTRL-O-SD, we finetune a pretrained Stable Diffusion model initialized from the stabilityai/stable-diffusion-2-1 checkpoint. As illustrated in Figure 3(a), CTRL-O extracts slots from a given image based on user-provided queries. These extracted slots are then incorporated into the caption, which is fed into Stable Diffusion. Notably, CTRL-O remains frozen during the fine-tuning process, distinguishing our approach from prior works like Slot Diffusion [44] and Stable LSD [19], where the object-centric model and the diffusion model are trained jointly.

Implementation Details We train the model on the COCO 2017 training set. For each image, we first extract COCO categories from its associated caption and use these categories to query CTRL-O, to generate the corresponding slots. These slots are subsequently appended to the caption, as shown in Figure 3(a). The resulting caption is then passed through the CLIP language encoder to condition Stable Diffusion. To integrate CTRL-O outputs into the CLIP language embedding space, we introduce a learnable linear layer that maps the extracted slots to the CLIP embedding space. During training, the only components updated are the

U-Net parameters of Stable Diffusion. We use random flips as a data augmentation strategy. Training is performed for 300 epochs with a learning rate of 2×10^{-5} , utilizing gradient accumulation with 2 steps. Additionally, we reproduce Stable LSD using the author-provided code and hyperparameters on the COCO dataset. The input resolution to the vision encoder for CTRL-O is 224×224 , while Stable LSD uses 448×448 .

H.2. Image Generation Metrics

Fréchet Inception Distance (FID) score We calculate the Fréchet Inception Distance (FID) score [17] to assess the quality of generated images in comparison to real images. The FID score computes the Fréchet distance between feature distributions of generated and real images, extracted via an Inception v3 model. Lower FID scores indicate a closer match to real images, corresponding to higher image fidelity and diversity.

CLIP-I Score We use CLIP-I Score to verify whether the generated images contain the same instances present in the query image. This should be the expected behavior of CTRL-O-SD when conditioned on a caption containing slots corresponding to specific instances. We compute this metric on the COCO validation set. We embed the generated image and the query image into the CLIP embedding space using the CLIP ViT Encoder (openai/clip-vit-base-patch16). We then compute the cosine similarity between the two embeddings. This similarity is averaged across all images to compute the final CLIP-I Score.

H.3. Image Reconstruction Visualization

In this section, we present a qualitative analysis of the reconstruction capabilities of the LSD and CTRL-O-SD models. The goal is to evaluate how effectively these models retain structural and semantic details. LSD provides the full 7-slot representation derived from the object-centric model to the generative model, providing comprehensive image context for reconstruction. In contrast, CTRL-O-SD provides the caption along with only a subset of slots corresponding to the categories in the caption to the generative model. To obtain these slots, we condition the slots in CTRL-O with the categories present in the caption and append the corresponding slots to the caption. This flexibility in CTRL-O-SD enables instance-specific image generation (see Fig. 4 for examples), which is not feasible with Stable LSD. As illustrated in Fig. 8, both models demonstrate comparable reconstruction quality.

H.4. Image Generation Failures

In this section we highlight some failure cases of CTRL-O-SD.

- **Incorrect Focus:** The model occasionally fails to accurately prioritize the main objects in the query, often diverting attention to irrelevant elements. For instance, when prompted to generate an image centered around a cell phone, the model might emphasize a person in the background instead. As we have seen from Table 1, CTRL-O does not achieve perfect binding. Hence, this failure can be caused by the slots not binding to the correct regions in the image.
- **Deformed Outputs:** The model sometimes generates distorted representations of people and animals, with unnatural proportions or malformed features. Such deformities highlight limitations in the model’s ability to represent detailed anatomy accurately, indicating a need for refined control over complex shapes and structures. This failure may also be attributed to the failures of the underlying generative model rather than CTRL-O.
- **Object Duplication:** There are instances where the model replicates objects within a single scene, leading to unrealistic and cluttered outputs.

These failure modes suggest areas for further improvement for CTRL-O and CTRL-O-SD, particularly in object binding for CTRL-O and image generation quality for CTRL-O-SD.

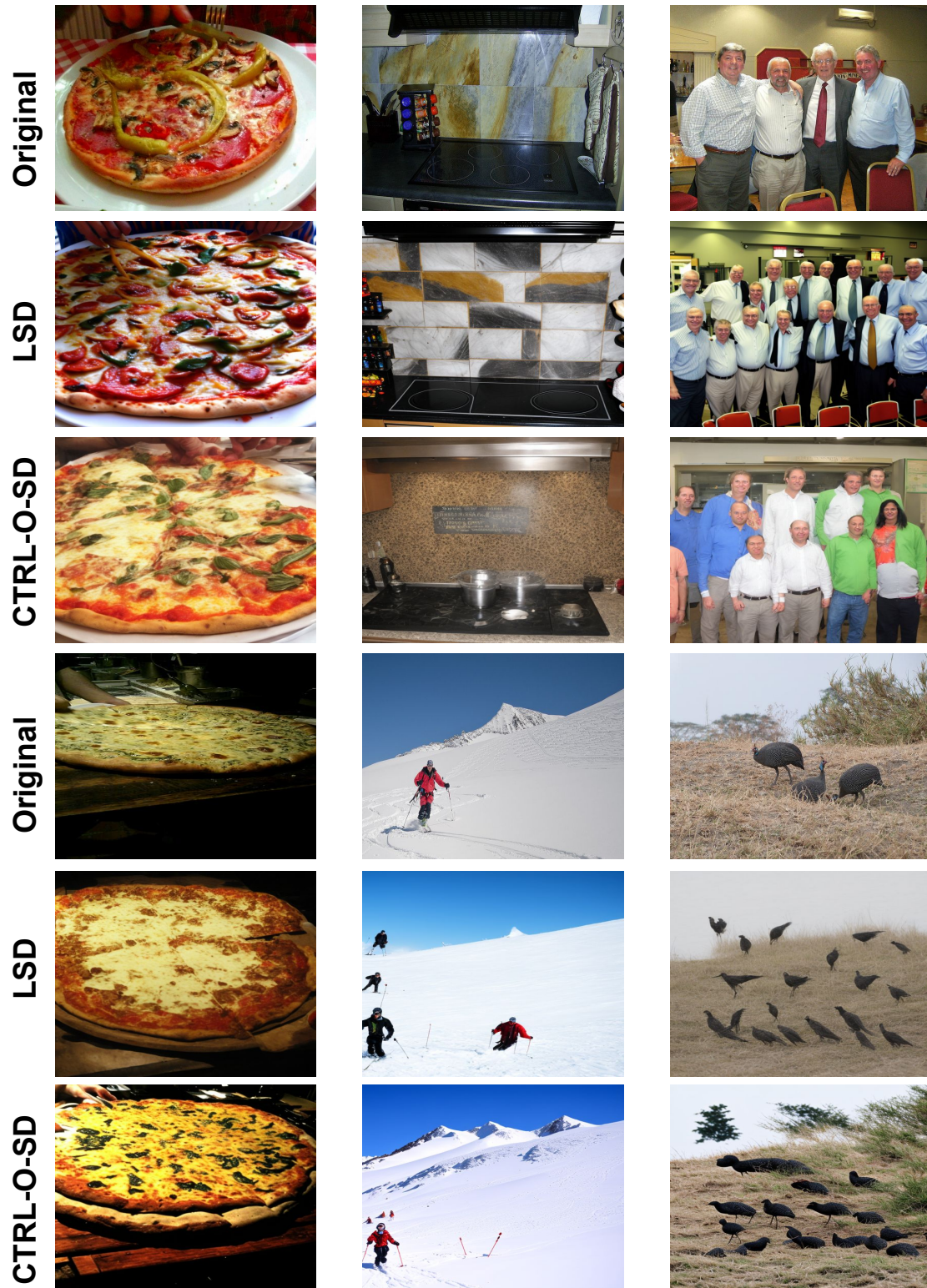


Figure 8. **Image Reconstruction** Qualitative comparisons of reconstruction outputs for LSD and CTRL-O-SD models. Each column corresponds to a different image. Rows correspond to original inputs, LSD generations, and CTRL-O-SD generations respectively. LSD generates outputs conditioned on full 7-slot representations derived from the original image, while CTRL-O-SD uses captions appended with a subset of slots for conditioning. The results show that both models achieve similar reconstruction quality.

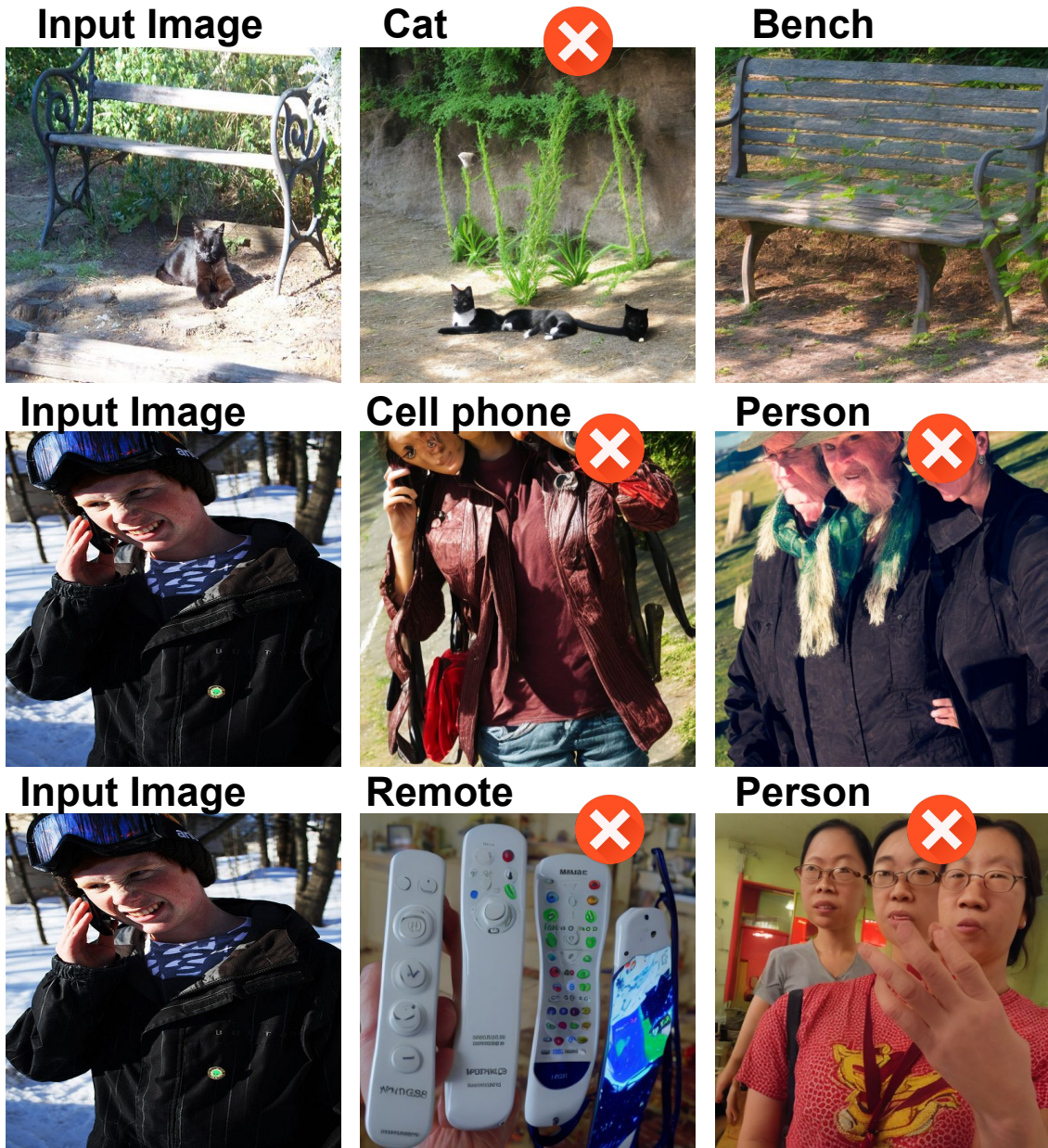


Figure 9. **Failure Modes of CTRL-O SD.** Examples highlighting some failures in CTRL-O-SD such as incorrect focus, deformities in representations of people or animals, and object duplication. Each labeled box illustrates specific instances of these failures.