

000
001
002
003
004
005
006
007
008
009
010
011054
055
056
057
058
059
060
061
062
063
064
065

Abstract

This paper studies the problem of object discovery – separating objects from the background without manual labels. Existing approaches rely on appearance cues, such as color, texture and location, to group pixels into object-like regions. However, by relying on appearance alone, these methods fail to reliably separate objects from the background in cluttered scenes. This is a fundamental limitation, since the definition of an object is inherently ambiguous and context-dependent. To resolve this ambiguity, in this work we choose to focus on dynamic objects – entities that are capable of moving independently in the world. We then scale the recent auto-encoder based frameworks for unsupervised object discovery from toy, synthetic images to complex, real world scenes by simplifying their architecture, and augmenting the resulting model with a weak learning signal from a motion segmentation algorithm. We demonstrate that, despite only capturing a small subset of the objects, this signal is enough to bias the model, which then learns to segment both moving and static instances of dynamic objects. We show that this model scales to our newly collected, photo-realistic synthetic dataset with street driving scenarios. Additionally, we leverage ground truth segmentation and flow annotations in this dataset for thorough ablation and evaluation. Finally, our experiments on the real-world KITTI dataset demonstrate that the proposed approach outperforms both heuristic- and learning-based methods by capitalizing on motion cues.

1. Introduction

Objects are the key building blocks of perception [31, 50]. We understand the world not in terms of pixels, surfaces, or entire scenes, but rather in terms of individual objects and their combinations. Object-centric representation makes tractable higher-level cognitive abilities such as causal reasoning, planning, etc., and are crucial for generalization and adaptation [5, 60]. In computer vision, progress has been achieved in object recognition recently [9, 24, 46], but these approaches rely on large amounts of expensive manual labels, and only cover a fixed vocabulary of object categories.

Anonymous CVPR submission

Paper ID 3978

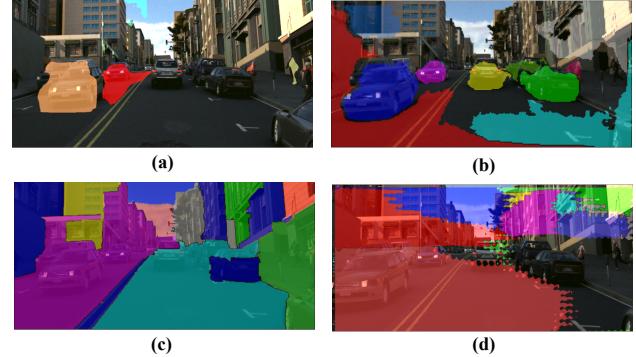


Figure 1. A sample from the PD dataset with: (a) motion segmentation from [14], top-10 segments produced by (b) our approach, (c) an heuristic-based MCG [3], and (d) learning-based SlotAttention [38]. Our method uses noisy, sparse motion segmentation to learn to separate both moving and static instances of dynamic objects from the background, whereas others cannot resolve the object definition ambiguity based on appearance alone.

Discovering objects and their extent in data – in a manner that generalizes across domains – remains an open problem.

What makes this task especially challenging is that the notion of an object is inherently ambiguous and context-dependent. Consider a car in Figure 1: its left door and the handle on that door can be treated as individual objects, or parts of the whole. It is thus not surprising that existing approaches that attempt to automatically separate objects from the background based on appearance struggle beyond controlled scenarios. In particular, classical methods that use graph-based inference tend to over- or under-segment the objects [3, 18] (Figure 1, bottom left). More recent learning-based methods model object discovery with structured generative networks, often leveraging iterative inference in the bottleneck of an auto-encoder [8, 16, 22, 37, 38]. While promising results have been demonstrated, these approaches are typically restricted to toy images with coloured geometric shapes on a plain background, and completely fail on realistic scenes (Figure 1, bottom right).

We posit that while the ambiguity of the object definition is not resolvable in the static image world without direct supervision, it has a natural resolution in the dynamic world

of videos. Concretely, we choose to focus on *dynamic* objects, which we define as entities that are capable of moving independently in the world. Independent object motion is a strong grouping cue, which has been shown to drive object learning in animal perception [13, 49]. In computer vision, there exists a long line of works on motion segmentation that automatically separate moving objects from the background based on optical flow [7, 14, 33, 41, 41, 61]. These methods have found numerous application in unsupervised [2, 43] and weakly-supervised machine learning algorithms [27, 44, 56].

In this work, we show how motion segmentation can be bootstrapped to group instances even when they are static. We build our approach on top of the framework for unsupervised object discovery proposed by Locatello et al. [38], and show how to scale it from toy images to realistic videos by leveraging independent object motion. We extend the architecture to videos of arbitrary length by introducing a spatio-temporal memory module [4], and simplify the grouping mechanism to scale the model to realistic scenes with large resolution and dozens of objects. We then demonstrate the importance of inductive biases based on independent object motion on the emergent representation and the extent to which it captures objects. In particular, we show how motion segments (Figure 1, top left) can guide the attention operation to discover static objects. Crucially, we show that motion segmentation of varying quality – even when sparse and noisy – can be sufficient to bias the model towards segmenting *both moving and static instances* (Figure 1, top right). Our approach only requires videos for training, and can segment objects in static images at inference time.

To go beyond the toy data used in [38], while still being able to thoroughly analyze the various aspects of the method, we have collected a new, photo-realistic, synthetic dataset using the ParallelDomain (PD) platform [1]. It consists of hundreds of videos, with crowded, street driving scenes, and comes with a full set of ground truth annotations, including object segmentations, 3D coordinates and optical flow, allowing us to ablate the importance of the quality of the motion segmentation to the final performance of the method. Finally, we demonstrate that the resulting method generalizes to real videos on the challenging KITTI dataset [19], and compare it to existing heuristic- and learning-based approaches.

2. Related work

In this work we study the problem of *object discovery* in realistic videos capitalizing on *motion segmentation* as a *learning signal for bottom-up grouping*. Below, we review the most relevant works in each of these areas.

Object discovery is the problem of separating objects from the background without manual labels. Traditional computer vision approaches treated it as perceptual grouping [36] – the idea that low and mid-level regularities in the data such as color, orientation and texture allow for approximately

parsing a scene into object-like regions. Notable approaches include [18], which uses graph-based inference to identify region boundaries, and [3] which first extracts regions on multiple scales with a normalized cut algorithm, and then groups them into object candidates. However, being purely appearance-based, these methods are not well equipped to resolve the inherent ambiguity of the object definition.

This problem has received renewed attention recently with the introduction of learning-based methods for object discovery [8, 16, 17, 22, 23, 29, 37, 38, 59]. A common approach is to use iterative inference to bind a set of variables to objects in an image [16, 22, 38], usually in a variational auto-encoder framework [35, 47]. A more efficient variant is proposed by Locatello et al. [38] in their SlotAttention framework. Concretely, they perform a single step of image encoding with a CNN (convolutional neural network) followed by an iterative attention operation, which is used to bind a set of variables, called slots, to image locations. The slots are then decoded individually and combined to reconstruct the image.

Many of the approaches above are capable of discovering objects in toy, synthetic scenes, but as we demonstrate in Section 4.5, they fail in more realistic environments, where appearance alone is not sufficient to separate the objects from the background. In this work, we extend SlotAttention to realistic videos by modifying the architecture of the model to allow it to scale to large scenes with dozens of objects, and incorporating inductive biases in the form of motion segmentation which naturally resolve the object/background ambiguity. Crucially, our method only uses motion segmentation as a sparse learning signal and the trained model does not rely on the motion information, being able to segment both moving and static instances.

Finally, several works have recently explored integrating inductive biases in the form of 3D geometry constraints [11, 15, 26, 51]. However, these methods remain limited to toy, synthetic environments. In contrast, our method uses independent object motion as a learning signal, allowing it to generalize to real-world scenes. Geometric priors are orthogonal to our approach and combining different forms of inductive biases is a promising direction for future work.

Motion segmentation is concerned with separating objects from the background using optical flow [28, 53, 55]. Early approaches [7, 33, 41, 41] tracked individual pixels with flow and then clustered the resulting trajectories inspired by the common fate principle [36]. While these methods have shown promising results on motion segmentation benchmarks, they do not generalize well in the wild due to their heuristic-based nature. More recently, several learning-based methods have been proposed [14, 61]. In particular, Dave et al. re-purpose a state-of-the-art object detection architecture [24] to detect and segment moving objects in an optical flow field. The model is trained on a toy, synthetic FlyingTh-

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

216 ings3D dataset [39], but can generalize to real videos due
 217 to appearance abstraction provided by the flow. We use this
 218 method in our work due to its high performance and simplicity
 219 combined with minimal supervision requirements. Note
 220 that since our method requires instance-level moving object
 221 masks, binary motion segmentation techniques [42, 57, 62]
 222 are not applicable in our scenario.
 223

224 **Learning from motion** is a paradigm inspired by evidence
 225 from cognitive science research, that independent object motion
 226 is a crucial cue for the development of the human visual
 227 system [49]. In computer vision, it has been adopted for
 228 weakly-supervised object detection [44] and semantic seg-
 229 mentation [27, 56], as well as for unsupervised representation
 230 learning [2, 43]. However, none of these works address the
 231 problem of object discovery from unlabeled videos. Yang
 232 et al. [63] use binary motion segmentation to train saliency
 233 models, but do not segment individual objects in complex
 234 scenes. Very recently, Tangemann et al. [54] have proposed
 235 to use motion segmentation to build compositional, genera-
 236 tive scene models. However, their approach employs motion
 237 segmentation as a pre-processing step during training and is
 238 not capable of object discovery at inference time.
 239

3. Method

240 In this section, we first introduce the SlotAttention frame-
 241 work for unsupervised object discovery, which serves as a
 242 basis for our approach, in Section 3.1. We then describe
 243 how we scale this architecture to real-world videos with
 244 dozens of objects in Section 3.2, and present our approach
 245 to incorporating independent motion priors in Section 3.3.
 246

3.1. Background

247 Following prior work [8, 22], SlotAttention [38] models
 248 object discovery as inference in an auto-encoder framework.
 249 Concretely, given an image $I \in \mathbb{R}^{H \times W \times 3}$, it is first passed
 250 through an encoder CNN to obtain a hidden representation
 251 $H = f_{enc}(I) \in \mathbb{R}^{H' \times W' \times D_{inp}}$. It is then processed by
 252 the attention module, which we describe below, to map H
 253 to a set of feature vectors of a fixed length K called slots
 254 $S \in \mathbb{R}^{K \times D_{slot}}$. Each slot S_i is broadcasted onto a 2D
 255 grid, and decoded individually with a decoder CNN $O_i =$
 256 $f_{dec}(S_i) \in \mathbb{R}^{H \times W \times 4}$, where the 4th dimension of the output
 257 represents the alpha mask A_i . Denoting the first 3 channels
 258 of O_i with I'_i , the complete image reconstruction is obtained
 259 via $I' = \sum_i A_i * I'_i$ and is used to supervise the model with
 260 an MSE (mean squared error) loss.
 261

262 The attention module is the key component of the ap-
 263 proach. It uses an iterative attention mechanism, similar to
 264 the one used in Transformer [58], to map from the input H to
 265 the slots S . In particular, the attention weights are computed
 266 with a dot product between the input features and slot states
 267 $W = \frac{1}{\sqrt{D}} k(H) \cdot q(S) \in \mathbb{R}^{N \times K}$, where k and q are learnable
 268

269 linear transformations and N is the spatial dimension of H .
 270 These attention weights are then normalized and used to com-
 271 pute the update values via $U = W^T v(H) \in \mathbb{R}^{K \times D}$, where
 272 W are the normalized attention weights, and v is another
 273 learnable linear transformation. A key difference to the clas-
 274 tical Transformer architecture is that the slots are initialized
 275 at random, and the inference is iterative. In particular, at ev-
 276 ery step l the slots are updated via $S_l = \text{update}(S_{l-1}, U_l)$,
 277 where the update function is implemented as a GRU [12]
 278 (gated recurrent unit).
 279

280 The intuition behind this approach is that the slots serve
 281 as a representational bottleneck and individual decoding of
 282 the slots results in them binding to spatially coherent regions,
 283 such as objects. Next, we describe how we modify the
 284 SlotAttention framework to scale it to real-world videos.
 285

3.2. A framework for object discovery in videos

286 Our model, shown in Figure 2, takes a sequence of
 287 video frames $\{I^1, I^2, \dots, I^T\}$ as input. Each frame is then
 288 processed by an encoder CNN, shown in yellow, to ob-
 289 tain an individual frame representation $H_t = f_{enc}(I^t) \in$
 290 $\mathbb{R}^{H' \times W' \times D_{inp}}$. These individual representations are ag-
 291 gregated by a ConvGRU spatio-temporal memory module [4]
 292 to obtain video encoding via $H'^t = \text{ConvGRU}(R^{t-1}, H^t)$,
 293 where $R^{t-1} \in \mathbb{R}^{H' \times W' \times D_{inp}}$ is the recurrent memory state.
 294

295 Next, we proceed to map the video representation H'^t to
 296 the set of slots S^t . It is easy to see, however, that the recur-
 297 rent slot assignment strategy proposed in [38] does not scale
 298 well to sequential inputs. Indeed, given a sequence of length
 299 T and L inference steps for each frame, the overall number
 300 of attention operations required to process the sequence is
 301 $T \times L$. Such a nested recurrence is both computationally
 302 inefficient, and can exacerbate the vanishing gradient prob-
 303 lem. To address this issue, as shown in the blue block in
 304 Figure 2, we only perform a single attention operation to
 305 directly compute the slot state $S^t = W^{t^T} v(H'^t) \in \mathbb{R}^{K \times D}$,
 306 where the attention matrix W^t is computed using the slot
 307 state in the previous frame S^{t-1} . For the first frame we use
 308 a learnable initial state S^0 .
 309

310 It is worth noting that the authors of [38] claim that it-
 311 erative inference on randomly initialized slots is crucial for
 312 the model to be able to generalize to a different number of
 313 objects at test time. However, we have found that simply
 314 increasing the number of slots to the maximum expected
 315 number of objects is sufficient to generalize to scenes of
 316 varying complexity. In that regard, our approach is similar
 317 to DETR [9], which also uses transformer query vectors as
 318 learnable object proposals that are capable of parsing both
 319 densely and sparsely populated scenes, but is trained in a
 320 fully supervised way.
 321

322 Finally, the resulting slot states S^t are processed with
 323 the decoder CNN, shown in green in Figure 2, to obtain the
 324 frame reconstruction. However, the individual slot decoding
 325

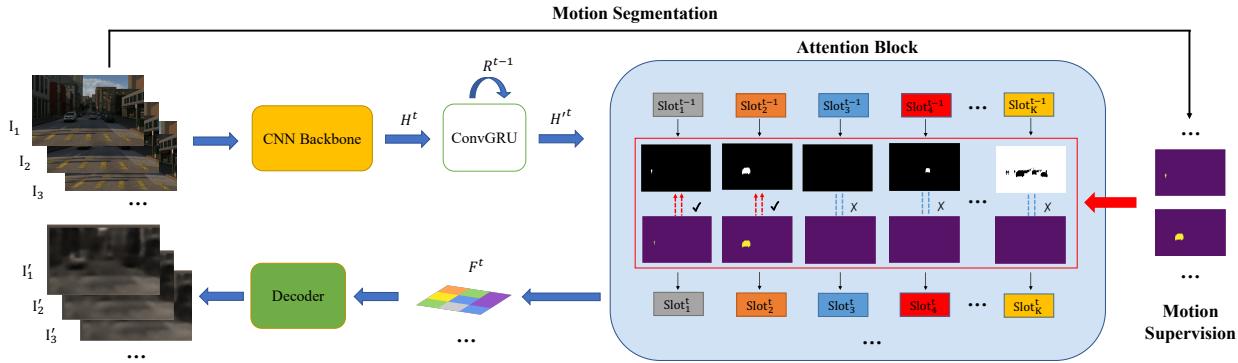
324
325
326
327
328
329
330
331
332
333
334
335
336378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Figure 2. Our method takes a sequence of frames as input and processes them individually with a backbone network (shown in yellow), and a ConvGRU recurrent memory module. The resulting feature maps H'^t are passed to the attention module (shown in blue) which binds them to a fixed set of slot variables via an attention operation. We additionally use automatically estimated motion segmentation to guide the attention operation for a subset of the slots. Finally, the slot states are combined in a single feature map F^t and decoded to reconstruct the frame. The reconstruction objective enforces generalization from moving to static instances.

approach from [38] does not scale well with the number of slots. Indeed, a full image reconstruction needs to be computed for each slot which quickly becomes prohibitively expensive in terms of memory, especially for large resolution frames. Instead, we propose to invert the order of slot decoding and slot recombination steps. In particular, we first broadcast each individual slot feature $S_i^t \in \mathbb{R}_{slot}^D$ to a feature map $F_i^t \in \mathbb{R}^{H' \times W' \times D_{slot}}$ and use the attention mask $W_{:,i}^t$ of the slot as an alpha mask A_i^t . We then construct a single output feature map $F^t = \sum_i A_i^t * F_i^t$, shown with a checkerboard pattern in the figure, and decode it via $I^t = f_{dec}(F^t) \in \mathbb{R}^{H \times W \times 3}$.

As we demonstrate in Section 4.3, the proposed single shot decoding strategy reduce the strength of the spatial cohesion prior in the original SlotAttention architecture, decreasing its object discovery capabilities. However, we also demonstrate that this naive prior does not generalize beyond toy, synthetic scenes. Instead, in the next section we describe our approach of incorporating an independent motion prior which provides a much stronger learning signal and works well with our efficient single shot decoding strategy.

3.3. Incorporating independent motion priors

Our method assumes that a set of sparse, instance-level motion segmentation masks $\mathcal{M} = \{M^1, M^2, \dots, M^{C^t}\}$ is provided with every video, with $M^t = \{m_1, m_2, \dots, m_{C^t}\}$, where C^t is the number of moving objects that were successfully segmented in frame t , and $m_j \in \{0, 1\}^{H' \times W}$ is a binary segmentation mask. Note that for every frame it is possible that $M_i = \emptyset$. This reflects the realistic assumption that a variable number of objects can be moving in any given frame and that in some frames all the objects can be static.

We propose to use these motion segmentation masks to directly supervise the slot attention maps $W^t \in \mathbb{R}^{N \times K}$, where $N = H' \times W'$. We thus need to map a variable number

of motion segmentations C^t to a fixed number of slots K in every frame. Following prior work on set-based supervision [9, 52], we first find an optimal bipartite matching between predicted and motion masks, and then optimize an object-specific segmentation loss. Specifically, we consider M^t also as a set of length K padded with \emptyset (no object). To find a bipartite matching between these two sets we search for a permutation of K elements with the lowest cost:

$$\hat{\sigma} = \arg \min_{\sigma} \sum_{i=1}^K \mathcal{L}_{seg}(m_i^t, W_{:, \sigma(i)}^t), \quad (1)$$

where $\mathcal{L}_{seg}(m_i^t, W_{:, \sigma(i)}^t)$ is the segmentation loss between the motion mask m_i^t and the attention map of the slot with index $\sigma(i)$. In practice, we efficiently approximate the optimal assignment with a greedy matching algorithm.

Once the assignment $\hat{\sigma}$ has been computed, the final motion supervision objective is defined as follows:

$$\mathcal{L}_{motion} = \sum_{i=1}^K \mathbb{1}_{\{m_i^t \neq \emptyset\}} \mathcal{L}_{seg}(m_i^t, W_{:, \hat{\sigma}(i)}^t). \quad (2)$$

That is, the loss is only computed for the slots for which a motion mask has been assigned, and the remaining slots are not constrained and can bind to any regions in the image. This is illustrated in the right part of Figure 2, where motion segmentation masks are available for only two objects in a crowded outdoor scene, and they get matched to the slots whose attention maps are most similar to the masks. Remaining slots are unconstrained, but still manage to capture both moving and static objects, as well as the background, driven by the image reconstruction objective. The actual segmentation loss \mathcal{L}_{seg} in Eq. 2 is the binary cross entropy:

$$\mathcal{L}_{seg}(m, W) = \sum_j -m_j \log(W_j) - (1-m_j) \log(1-W_j), \quad (3)$$

432 where N is the spatial dimension of the attention map W .
 433

434 3.4. Loss function and optimization

435 Our final objective is defined as follows:
 436

$$437 \quad \mathcal{L} = \mathcal{L}_{recon} + \lambda_M \mathcal{L}_{motion} + \lambda_T \mathcal{L}_{temp}, \quad (4)$$

438 where \mathcal{L}_{recon} is the MSE loss for the image reconstruction,
 439 \mathcal{L}_{temp} is a temporal consistency regularization term, and
 440 λ_M and λ_T are the weights for the motion supervision and
 441 temporal consistency terms. The latter is defined as
 442

$$443 \quad \mathcal{L}_{temp}(S) = \sum_{t=1}^{T-1} \|\mathbb{I} - \text{softmax}(S^t \cdot (S^{t+1})^\mathbf{T})\|, \quad (5)$$

444 where $\mathbb{I} \in \mathbb{R}^{K \times K}$ is the identity matrix. It is easy to see that
 445 this term is a form of a temporal contrastive loss encouraging
 446 similarity between feature representations of the slots in
 447 consecutive frames and thus improving the temporal consistency
 448 of the slot bindings. The model is trained on video clips of
 449 length T and we ensure that at least half of the clips in a
 450 batch have a non-empty set of motion segmentations \mathcal{M} .
 451

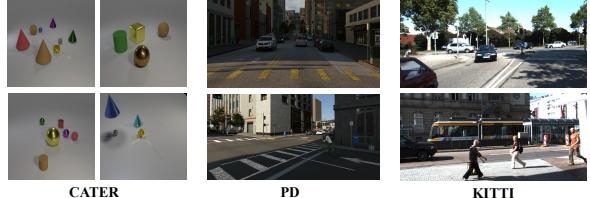
452 4. Experimental evaluation

453 4.1. Datasets and evaluation

454 We use two synthetic datasets for the analysis of the pro-
 455 posed approach: CATER [20] for ablating the architecture of
 456 the model and a realistic ParallelDomain (PD) dataset which
 457 we collect ourselves and use to analyze the impact of the
 458 motion segmentation quality on the model’s performance.
 459 In addition, we use a real-world KITTI benchmark [19] for
 460 comparison to the state of the art.
 461

462 **CATER** is a video version of the CLEVR [30] dataset which
 463 was used in many recent works on unsupervised object dis-
 464 covery [8, 29, 38]. We used the provided engine to generate
 465 2,000 videos by placing between 4 and 8 geometric shapes,
 466 such as cubes or cones, on a plain background at random,
 467 and assigning a random color to each instance. Each object
 468 can then move on a random trajectory or remain static, and
 469 the camera motion is also randomized. We use 1600 videos
 470 for training and 400 for evaluation, with each video being
 471 40-frames long with a resolution 128×128 . Dataset exam-
 472 ples are shown in Figure 3, left. For ablation analysis, we
 473 randomly assign one object as moving in each video and
 474 use the ground truth mask of that object as a motion mask.
 475 Notice that we do experiment with automatically estimated
 476 motion segmentation on more challenging PD and KITTI.
 477

478 **ParallelDomain (PD)** is our synthetic dataset with street
 479 driving scenarios (see Figure 3, center). It was collected
 480 using a state-of-the-art synthetic data generation service [1].
 481 The training set contains 154 photo-realistic videos with
 482 driving scenarios in city environments captured at 20 FPS.
 483



484 Figure 3. Frame samples from the video datasets used in our
 485 experiments. CATER [20] (left) is a toy, synthetic dataset similar
 486 to the ones used in prior works. PD (center) is our own collection
 487 of photo-realistic, synthetic videos, which is a major step forward
 488 in visual complexity. KITTI [19] (right) is a real world benchmark
 489 with outdoor scenes.
 490

491 Each video is 10 seconds long and comes with 6 independent
 492 camera views, effectively increasing the dataset size to 924
 493 videos. We use 51 videos from a disjoint set of scenes for
 494 evaluation. Each video comes with a full set of ground truth
 495 annotations, including optical flow, allowing us to conduct
 496 a detailed analysis of the impact of the motion segmen-
 497 tation quality on our method’s performance. More statistics
 498 and qualitative examples are provided in the supplemen-
 499 tary material. The dataset will be released with the paper.
 500

501 **KITTI** is a real-world benchmark with city driving scenarios
 502 which comes with a variety of annotations (Figure 3, right).
 503 In this work, we use the instance segmentation subset of the
 504 dataset for evaluation, which consists of 200 frames, which
 505 we resize to 368×1248 . Notice that instance segmentation
 506 annotations are provided on individual images in this dataset,
 507 without the temporal context, allowing us to demonstrate
 508 that our approach does not require videos at inference time.
 509 All the categories labeled in KITTI correspond to dynamic
 510 objects, so we use the whole set of annotations for evalua-
 511 tion. Since our approach does not require any labels for training,
 512 we use all the 147 videos in the training set of KITTI to learn
 513 to discover the objects in the real world.
 514

515 **Evaluation metrics.** We use Adjusted Rand Index (ARI)
 516 as the main metric for comparing object discovery capa-
 517 bilities of the models. ARI is a clustering similarity met-
 518 ric which captures how well predicted segmentation masks
 519 match ground-truth masks in a permutation-invariant fash-
 520 ion. This is more suitable for evaluation of unsupervised
 521 approaches than, say, mIoU, because it does not require for
 522 the methods to make the decision which segments represent
 523 the objects and which correspond to the background. Fol-
 524 lowing prior work [22, 38], we only measure ARI based on
 525 foreground objects, which we refer to as Fg. ARI.
 526

527 4.2. Implementation details

528 For the components of our model shared with SlotAtten-
 529 tion [38] we follow their architecture and training protocol
 530 exactly, and describe the remaining details below. We will
 531 release all the code, models, and synthetic data to ensure
 532

540 reproducibility of the experiments.
 541

542 We replace the shallow encoder used in [38] with a
 543 ResNet18 [25] to scale the representational power to re-
 544 alistic scenes. We also experiment with deeper backbones
 545 in the supplementary material. All the models are trained
 546 from scratch unless stated otherwise. We additionally report
 547 results with contrastive-learning pre-training in the supple-
 548 mentary. To be able to capture small objects, we remove the
 549 last 2 max pooling layers from the ResNet, and add corre-
 550 sponding dilation ratio to preserve the field of view. We use
 551 10 slots for the experiments on CATER and 45 slots on PD
 552 and KITTI to account for the larger number of objects.
 553

554 All the models are trained for 500 epochs using
 555 Adam [34] with a batch size 20 and learning rate 0.001.
 556 Following [38], we use learning rate warm-up [21] and an
 557 exponential decay schedule to prevent early saturation and
 558 reduce variance. We set λ_M to 0.5 and λ_T to 0.01 on the
 559 validation set of CATER, and use these value in all the ex-
 560 periments. The video-based variants are trained using clips
 561 of length 5. At inference time, the model is evaluated in a
 562 sliding window fashion with a stride 5.
 563

564 We experiment with two motion segmentation algorithms
 565 – an heuristic-based [33], and a learning-based one [14], for
 566 which we only use the motion stream trained on the toy Fly-
 567 ingThings3D dataset [39]. Both methods take optical flow
 568 as input, so we evaluate them with both ground truth flow,
 569 and flow estimated with the state-of-the-art supervised [55]
 570 and unsupervised [53] approaches. Since the outputs of
 571 both methods contain many noisy segments, we apply a few
 572 generic post-processing steps to clean up the results. They
 573 removing very large and very small segments, as well as
 574 segments at the image boundary. The details of the post-
 575 processing are provided in the supplementary material.
 576

577 We compare our approach to several recent learning-
 578 based object discovery approaches as well as to a classical,
 579 heuristic-based method. In particular, we choose SlotAtten-
 580 tion [38], MONet [18], SCALOR [29], and S-IODINE [22]
 581 as a representative sample of learning-based methods, with
 582 S-IODINE also being a video-based approach. For each
 583 of these methods, we replace the original backbone with
 584 ResNet18 and match the input resolution to the one used
 585 by our method for a fair comparison, but keep all the other
 586 details intact. All the models are trained until convergence.
 587 We use MCG [3] as an heuristic-based baseline. It is a pro-
 588 posal generation method, so to obtain a single interpretation
 589 of an image, we sample the top scoring proposals until all
 590 the pixels are covered. For overlapping segments, we assign
 591 the corresponding pixels to the smaller segment.
 592

593 4.3. Architectural analysis

594 In this section, we begin the analysis of our method by
 595 studying the variants of the auto-encoder framework for ob-
 596 ject discovery on the validation set of CATER in Table 1.
 597

ConvGRU	Slot inf.	Temp.	Decode.	Motion	Recon.	Fg. ARI	594
–	Iter.	X	Per slot	X	✓	64.4	595
frame	Iter.	X	Per slot	X	✓	66.3	596
clip	Iter.	X	Per slot	X	✓	71.5	597
clip	1-shot	X	Per slot	X	✓	83.2	598
clip	1-shot	✓	Per slot	X	✓	86.7	599
clip	1-shot	✓	1-shot	X	✓	34.5	600
clip	1-shot	✓	1-shot	✓	✓	92.7	601
clip	1-shot	✓	1-shot	✓	X	77.9	602

603 Table 1. Analysis of the model architecture using Fg. ARI on
 604 the validation set of CATER. We ablate the ConvGRU module,
 605 slot inference strategy, temporal consistency constraint, decoding
 606 strategy, independent motion prior, and the reconstruction objective.
 607 Combining motion priors with reconstruction leads to best results.
 608

609 Firstly, we evaluate the original SlotAttention model (row
 610 1 in the table), which serves as a basis for our approach,
 611 and find that it performs reasonably well on this toy dataset,
 612 though the Fg. ARI scores are noticeable lower than those
 613 reported in the original paper [38] on CLEVR. This is ex-
 614 plained by the fact that the scenes in CATER are more chal-
 615 lenging, with a larger variance in the number of objects and
 616 more occlusions.
 617

618 Next, we convert the frame-level architecture of SlotAt-
 619 tention to a video-level model by adding a ConvGRU after
 620 the encoder. This has only a minor effect on the performance
 621 when trained on 1-frame sequences (row 2 in the table), but
 622 training on video clips (row 3) results in a 5.2 points increase
 623 in Fg. ARI score. This demonstrates that the feature space
 624 of the recurrent model can capture video dynamics and thus
 625 simplify separating objects from the background.
 626

627 However, going from single frame inputs to clips in-
 628 creases the memory requirements of the model. To mitigate
 629 this issue, we now study the architectural modifications pro-
 630 posed in Section 3.2. Firstly, replacing iterative inference on
 631 randomly initialized slots with a single attention operation
 632 with a learnable initialization not only results in an improved
 633 computational efficiency, but also significantly improves the
 634 performance. Incorporating the temporal consistency term
 635 in the loss further boosts the Fg. ARI score due to more
 636 robust slot binding. Next, switching to 1-shot decoding sig-
 637 nificantly reduces the memory consumption of the model,
 638 but also results in it largely loosing its object-discovery cap-
 639 abilities. This demonstrates that individual slot decoding
 640 was crucial for enforcing the spatial cohesion prior of the
 641 SlotAttention model.
 642

643 Despite this disadvantage, incorporating a weak learn-
 644 ing signal in the form of a motion segmentation not only
 645 recovers, but significantly improves the model’s perfor-
 646 mance. This demonstrates that independent motion is a much
 647 stronger and more generic prior than appearance and loca-
 648 tion similarity used in the SlotAttention, even in a toy dataset
 649 like CATER. Finally, the last row of Table 1 shows that the
 650 reconstruction objective is still important for achieving top
 651 performance by enforcing generalization from moving to
 652 static instances.
 653

648	Model	Motion seg.	Fg. ARI Stat.	Fg. ARI Mov.	Fg. ARI All
649	Ours	None	10.5	18.4	13.1
650	Ours	GT all	68.0	72.5	71.7
651	Ours	GT moving	48.4	66.7	59.6
652	Ours	GT flow + [33]	36.9	47.5	42.8
653	Ours	GT flow + [14]	47.0	57.3	51.7
654	Ours	RAFT flow + [14]	45.6	56.7	50.9
655	Ours	SMURF flow + [14]	44.1	56.1	50.5
656	-	RAFT flow + [14]	2.7	5.3	3.4

Table 2. Analysis of the effect of the quality of motion segmentation on the model’s performance on the validation set of PD. We gradually reduce the quality of the motion segments starting from ground truth to fully estimated. Our method learns to discover both moving and static instances guided by a very sparse motion signal.

4.4. Object discovery in realistic videos

We now explore how well the model introduced above scales to realistic street-driving scenes in the PD dataset in Table 2 and Figure 4. We separately report the Fg. ARI score for moving and static objects to asses the network’s generalization abilities. We begin evaluating the baseline variant of our model without independent motion priors, and observe that appearance similarity is indeed not sufficient for object discovery in complex scenes with cluttered background, as reflected by the low Fg. ARI score. Qualitatively, first column of Figure 4 illustrates that this variant completely fails to discover any objects, and instead segments the scene into random patches based on color and location similarity.

Next, we establish the upper bound for our model’s performance by using all the ground truth object segments (corresponding to moving and static objects) for training. This fully-supervised approach reaches an Fg. ARI score of 71.7, which is significantly below the 92.7 obtained by the best version of our model on CATER, further emphasizing the complexity of PD. Qualitatively, as can be seen in the second column of Figure 4, this variant successfully captures all the clearly visible objects in a scene, and also groups the background pixels together.

Only using the ground truth segments corresponding to the moving objects, which simulates the theoretical scenario in which we have a perfect motion segmentation algorithm, does result in a performance drop of 11.3 Fg. ARI points, which is especially noticeable for static objects, but the overall score remains 46.5 points higher than the baselines trained without the motion prior. Qualitatively, the model is able to accurately segment most of the moving and static instances, as shown in the third column in Figure 4. However, this variant oversegments the background, demonstrating that explaining as many objects in the scene as possible is crucial for learning a strong background model.

Switching to actual motion segmentation algorithms, we first compare the state-of-the-art heuristic-based and learning-based methods using the ground truth optical flow as input in rows 5 and 6 of the Table 2. As expected, we observe that the more recent learning-based method produces

	Learning-based	PD	KITTI	702
SlotAttention [38]	✓	10.2	13.8	703
MONet [8]	✓	11.0	14.9	704
SCALOR [29]	✓	18.6	21.1	705
S-IODINE [22]	✓	9.8	14.4	706
MCG [3]	✗	25.1	40.9	707
Ours	✓	50.9	47.1	708

Table 3. Comparison to the state-of-the-art approaches for object discovery on the validation sets of PD and KITTI using Fg. ARI. Our approach outperforms both learning- and heuristic-based methods by capitalizing on independent motion cues.

more accurate motion segmentations, which in turn results in a higher performance of our approach. Qualitatively, this model, shown in the 4th column in Figure 4, has a slightly lower recall than the variant trained with ground truth moving segments due to the sparser learning signal. Intriguingly, replacing ground truth flow with the one estimated with a state-of-the-art supervised RAFT [55], or self-supervised SMURF [53] algorithms barely changes the performance, despite a noticeable decrease in the motion segmentation quality (last column in Figure 4). This result demonstrates the robustness of our method to noise. We use RAFT flow for the remainder of the experiments.

Finally, to better quantify the ability of our model to generalize from sparse, noisy motion segmentations to the whole distribution of objects in crowded scenes, we evaluate the Fg. ARI score of the motion segmentations themselves in the last row of Table 2. We can see that these masks indeed mostly capture the moving objects, however, even for those only a tiny fraction is segmented. In contrast, our approach, capitalizing on this noisy and incomplete signal, increases the overall ARI score by a factor of 15.

4.5. Comparison to the state of the art

Finally, we compare our approach to the state-of-the-art on the validation sets of PD and KITTI in Table 3. Firstly, we observe that all the learning-based methods fail to achieve non-trivial results on both datasets. This confirms our hypothesis that appearance alone is not a sufficient signal to separate objects from the background in realistic environments. In contrast, our proposed approach outperforms all these methods by a wide margin by capitalizing on independent motion cues.

Interestingly, the classical MCG approach performs significantly better than the more recent learning-based methods. Our method outperforms MCG on both datasets, with the margin being significantly larger on PD. Recall that KITTI is an image-based benchmark, where the annotated frames are selected to prominently feature the objects of interest. In contrast, PD is a densely labeled video dataset with more challenging camera angles and more background clutter (see Figure 5 for a qualitative comparison).

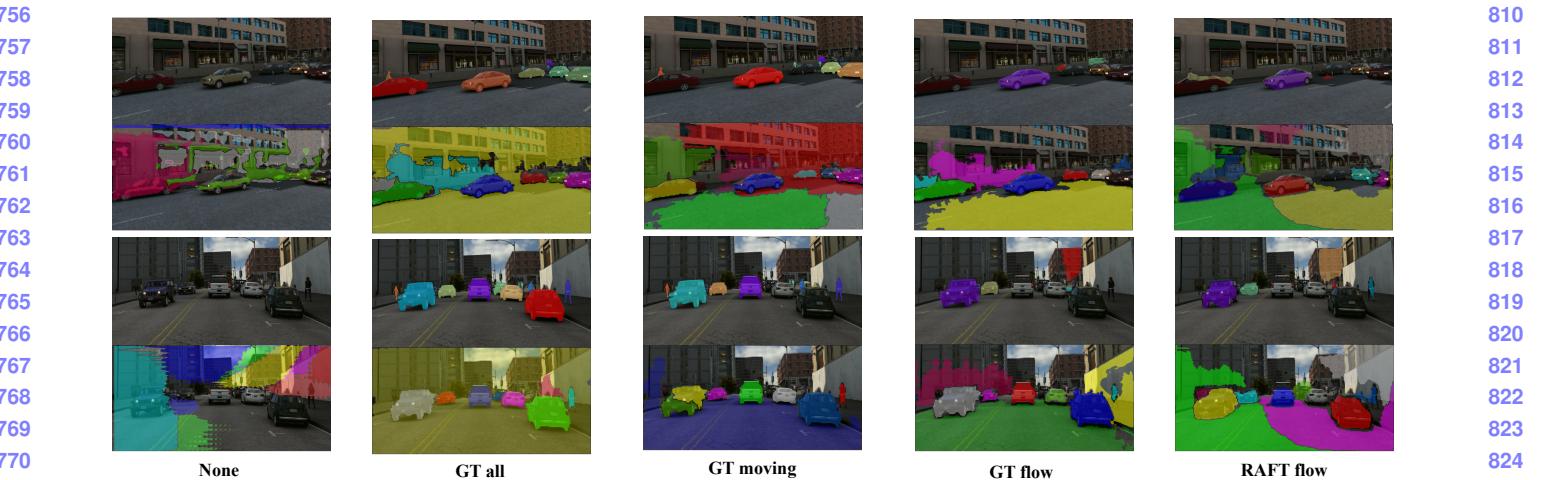


Figure 4. Top-10 masks produced by our model with varying quality of motion priors on the validation set of PD. We show the motion masks used for supervision on top of the corresponding model’s outputs. For the last two columns the approach of Dave et al. [14] is used to compute the motion segmentation.

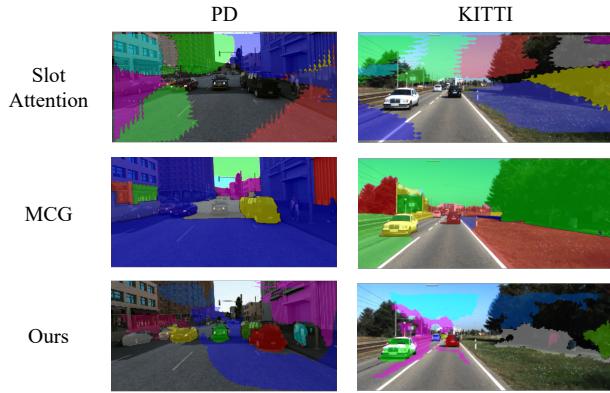


Figure 5. Qualitative comparison of our approach and representative heuristic- and learning-based methods on the validation sets of PD and KITTI (showing top-10 masks). Ours learns to successfully separates objects from the background, whereas appearance-based methods struggle in cluttered environments.

5. Discussion and limitations

Discovering objects and their extent from raw data is a challenging problem due to the ambiguity of what constitutes an object. In this work, we propose one way to automatically resolve this ambiguity by focusing on dynamic objects and using independent motion as an inductive bias in an auto-encoder framework. Our analysis demonstrates promising results in real-world environments, while further raising a number of important questions.

Generalization to non-dynamic objects. While independent object motion provides a convenient signal for object discovery from data, it ignores objects that are not capable of moving by themselves, but might be important for downstream tasks. In particular, in indoor environments people

interact with accessories, electronics, food, etc., and capturing these objects is crucial for action recognition [45, 65] and robotics [6, 40]. Notice, however, that extending the definition of a dynamic object to those entities that either move by themselves or can be moved by humans covers most of such cases. Classical motion segmentation approaches [7, 33] do attempt to capture all the objects that fall into this more general definition, but do not generalize in the wild. Developing more robust, learning-based versions of these methods is a critical step towards a generic object discovery algorithm.

Object category imbalance in the real world. Like any other learning-based method, ours is susceptible to focusing on the most common categories, while ignoring the objects in the tail of the distribution. For instance, in the real world we might see lots of moving people, vehicles and animals, and sometimes a person picking up a piece of litter. In theory, this should allow our method to discover not only what people, cars and animals are, but also litter. However, it might happen too infrequently in practice. Fortunately, this problem has received a lot of attention in the few-shot and continual learning domains [10, 32, 48, 64], and the proposed solutions can be integrated into our framework.

Supervision used to train the motion segmentation algorithm. The approach of Dave et al. [14], used in our experiments, is trained on the toy, synthetic FlyingThings3D [39] dataset with ground truth moving object masks. This raises the question of whether it is this indirect object-level supervision which makes our method outperform other, fully unsupervised approaches. To address this concern, in the supplementary material we directly pre-train SlotAttention on FlyingThings3D in a fully-supervised way, showing this does not have a significant effect on its object discovery performance in realistic videos due to the large domain gap.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

References

- [1] Parallel domain. <https://paralleldomain.com/>, November 2021. 2, 5
- [2] Pulkit Agrawal, Joao Carreira, and Jitendra Malik. Learning to see by moving. In *CVPR*, 2015. 2, 3
- [3] Pablo Arbeláez, Jordi Pont-Tuset, Jonathan T Barron, Ferran Marques, and Jitendra Malik. Multiscale combinatorial grouping. In *CVPR*, 2014. 1, 2, 6, 7
- [4] Nicolas Ballas, Li Yao, Chris Pal, and Aaron Courville. Delving deeper into convolutional networks for learning video representations. In *ICLR*, 2016. 2, 3
- [5] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019. 1
- [6] Aude Billard and Danica Kragic. Trends and challenges in robot manipulation. *Science*, 364(6446), 2019. 8
- [7] Thomas Brox and Jitendra Malik. Object segmentation by long term analysis of point trajectories. In *ECCV*, 2010. 2, 8
- [8] 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. 1, 2, 3, 5, 6, 7
- [9] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 1, 3, 4
- [10] Nadine Chang, Zhiding Yu, Yu-Xiong Wang, Anima Anandkumar, Sanja Fidler, and Jose M Alvarez. Image-level or object-level? A tale of two resampling strategies for long-tailed detection. In *ICML*, 2021. 8
- [11] Chang Chen, Fei Deng, and Sungjin Ahn. Object-centric representation and rendering of 3D scenes. *arXiv preprint arXiv:2006.06130*, 2020. 2
- [12] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *EMNLP*, 2014. 3
- [13] M. Cynader, N. Berman, and A. Hein. Cats reared in stroboscopic illumination: Effects on receptive fields in visual cortex. *Proceedings of the National Academy of Sciences*, 70(5):1353–1354, 1973. 2
- [14] Achal Dave, Pavel Tokmakov, and Deva Ramanan. Towards segmenting anything that moves. In *ICCV Workshops*, 2019. 1, 2, 6, 7, 8
- [15] Yilun Du, Kevin Smith, Tomer Ulman, Joshua Tenenbaum, and Jiajun Wu. Unsupervised discovery of 3D physical objects from video. In *ICLR*, 2021. 2
- [16] Martin Engelcke, Adam R Kosiorek, Oiwi Parker Jones, and Ingmar Posner. GENESIS: Generative scene inference and sampling with object-centric latent representations. In *ICLR*, 2020. 1, 2
- [17] SM Eslami, Nicolas Heess, Theophane Weber, Yuval Tassa, David Szepesvari, Geoffrey E Hinton, et al. Attend, infer,
- repeat: Fast scene understanding with generative models. In *NeurIPS*, 2016. 2
- [18] Pedro F Felzenszwalb and Daniel P Huttenlocher. Efficient graph-based image segmentation. *IJCV*, 2004. 1, 2
- [19] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *CVPR*, 2012. 2, 5
- [20] Rohit Girdhar and Deva Ramanan. CATER: A diagnostic dataset for compositional actions and temporal reasoning. In *ICLR*, 2020. 5
- [21] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large mini-batch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 6
- [22] 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. 1, 2, 3, 5, 6, 7
- [23] Klaus Greff, Antti Rasmus, Mathias Berglund, Tele Hao, Harri Valpola, and Jürgen Schmidhuber. Tagger: Deep unsupervised perceptual grouping. In *NeurIPS*, 2016. 2
- [24] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017. 1, 2
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *ECCV*, 2016. 6
- [26] Paul Henderson and Christoph H Lampert. Unsupervised object-centric video generation and decomposition in 3D. In *NeurIPS*, 2020. 2
- [27] Seunghoon Hong, Donghun Yeo, Suha Kwak, Honglak Lee, and Bohyung Han. Weakly supervised semantic segmentation using web-crawled videos. In *CVPR*, 2017. 2, 3
- [28] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, 2017. 2
- [29] Jindong Jiang, Sepehr Janghorbani, Gerard De Melo, and Sungjin Ahn. SCALOR: Generative world models with scalable object representations. In *ICLR*, 2020. 2, 5, 6, 7
- [30] 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*, 2017. 5
- [31] Daniel Kahneman, Anne Treisman, and Brian J Gibbs. The reviewing of object files: Object-specific integration of information. *Cognitive psychology*, 1992. 1
- [32] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *ICLR*, 2020. 8
- [33] Margret Keuper, Bjoern Andres, and Thomas Brox. Motion trajectory segmentation via minimum cost multicut. In *ICCV*, 2015. 2, 6, 7, 8
- [34] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 6

- 972 [35] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 2 1026
- 973 [36] Kurt Koffka. *Principles of Gestalt psychology*. Routledge, 1027
- 974 2013. 2 1028
- 975 [37] Zhixuan Lin, Yi-Fu Wu, Skand Vishwanath Peri, Weihao Sun, 1029
- 976 Gautam Singh, Fei Deng, Jindong Jiang, and Sungjin Ahn. 1030
- 977 Space: Unsupervised object-oriented scene representation via 1031
- 978 spatial attention and decomposition. In *ICLR*, 2020. 1, 2 1032
- 979 [38] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, 1033
- 980 Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, 1034
- 981 Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning 1035
- 982 with slot attention. In *NeurIPS*, 2020. 1, 2, 3, 4, 5, 6, 1036
- 983 7 1037
- 984 [39] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, 1038
- 985 Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A 1039
- 986 large dataset to train convolutional networks for disparity, 1040
- 987 optical flow, and scene flow estimation. In *CVPR*, 2016. 3, 6, 1041
- 988 8 1042
- 989 [40] Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof 1043
- 990 grasnet: Variational grasp generation for object manipulation. 1044
- 991 In *ICCV*, 2019. 8 1045
- 992 [41] Peter Ochs and Thomas Brox. Higher order motion models 1046
- 993 and spectral clustering. In *CVPR*, 2012. 2 1047
- 994 [42] Anestis Papazoglou and Vittorio Ferrari. Fast object 1048
- 995 segmentation in unconstrained video. In *ICCV*, 2013. 3 1049
- 996 [43] Deepak Pathak, Ross Girshick, Piotr Dollár, Trevor Darrell, 1050
- 997 and Bharath Hariharan. Learning features by watching objects 1051
- 998 move. In *CVPR*, 2017. 2, 3 1052
- 999 [44] Alessandro Prest, Christian Leistner, Javier Civera, Cordelia 1053
- 1000 Schmid, and Vittorio Ferrari. Learning object class detectors 1054
- 1001 from weakly annotated video. In *CVPR*, 2012. 2, 3 1055
- 1002 [45] Siyuan Qi, Wenguan Wang, Baoxiong Jia, Jianbing Shen, 1056
- 1003 and Song-Chun Zhu. Learning human-object interactions by 1057
- 1004 graph parsing neural networks. In *ECCV*, 2018. 8 1058
- 1005 [46] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 1059
- 1006 Faster R-CNN: Towards real-time object detection with region 1060
- 1007 proposal networks. *NeurIPS*, 2015. 1 1061
- 1008 [47] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. 1062
- 1009 Stochastic backpropagation and approximate inference 1063
- 1010 in deep generative models. In *ICML*, 2014. 2 1064
- 1011 [48] Dvir Samuel and Gal Chechik. Distributional robustness loss 1065
- 1012 for long-tail learning. In *ICCV*, 2021. 8 1066
- 1013 [49] Elizabeth S Spelke. Principles of object perception. *Cognitive 1067*
- 1014 science, 1990. 2, 3 1068
- 1015 [50] Elizabeth S Spelke and Katherine D Kinzler. Core knowledge. 1069
- 1016 *Developmental science*, 2007. 1 1070
- 1017 [51] Karl Stelzner, Kristian Kersting, and Adam R Kosiorek. 1071
- 1018 Decomposing 3D scenes into objects via unsupervised volume 1072
- 1019 segmentation. *arXiv preprint arXiv:2104.01148*, 2021. 2 1073
- 1020 [52] Russell Stewart, Mykhaylo Andriluka, and Andrew Y Ng. 1074
- 1021 End-to-end people detection in crowded scenes. In *CVPR*, 1075
- 1022 2016. 4 1076
- 1023 [53] Austin Stone, Daniel Maurer, Alper Ayvaci, Anelia Angelova, 1077
- 1024 and Rico Jonschkowski. SMURF: Self-teaching multi-frame 1078
- 1025 unsupervised RAFT with full-image warping. In *CVPR*, 2021. 2, 6, 7 1079