

000
001
002
003
004
005
006
007
008
009
010
011054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Track, Check, Repeat: An EM Approach to Unsupervised Tracking

Anonymous CVPR 2021 submission

Paper ID 7264

Abstract

We propose an unsupervised method for 3D object segmentation and motion/content disentanglement of unlabelled RGB-D video streams. We isolate the independently-moving foreground under an arbitrary moving camera, split it into trackable components, and track each component in 3D, across partial and full occlusions. The method works by iteratively learning 2D and 3D identity detectors, objectness detectors, bottom-up 2D and 3D motion estimators, and a motion prior; together these modules form a high-precision tracker. Learning happens in an expectation-maximization framework, where in the expectation step we fire all modules and look for agreement among them, and in the maximization step we re-train the modules to improve this agreement. This iterative process gradually expands recall, until the entire video is explained by tracklets. The constraint of ensemble agreement helps combat contamination of the pseudo-labels, and standard data augmentation techniques help the modules generalize to yet-unlabelled data. We compare against existing state-of-the-art unsupervised object discovery and tracking methods, using challenging real-world videos from the KITTI benchmark, and show strong improvements.

1. Introduction

Humans can detect moving objects and delineate their approximate extent in 3D from a single and possibly moving viewpoint, and without ever been supplied 3D boxes or 3D segmentation masks as supervision, whether in their lifespans or their evolutionary history. How does this remarkable ability develop without supervision? Neuroscience and psychology literature points to a variety of perceptual grouping cues, which makes some regions look more “object-like” than others. These types of objectness cues have long been known in computer vision literature; why has this domain knowledge not yet led to powerful self-supervised object detectors?

Work on integrating perceptual grouping cues into computer vision models stretches back decades, and is still

likely serves as inspiration for many of the design decisions in modern computer vision architectures related to attention, metric learning, two-stream architectures, and so on. Much of the current work on object recognition and tracking is fully-supervised, and relies on vast pools of human-provided annotations. On the unsupervised side, a variety of deep learning-based methods have been proposed, which hinge on reconstruction objectives and object-centric or scene-centric bottlenecks in the architecture [5]. These methods are rapidly advancing, but so far only on toy worlds, made up of simple 2D or 3D shapes against simple backgrounds – a far cry from the complexity tackled in older works [13].

Classic methods of object discovery, such as center-surround saliency in color or flow [1], are known to be brittle, but they need not be discarded entirely. We suggest to mine and exploit the admittedly rare success scenarios of these models, to bootstrap the learning of something more general. We hypothesize that if the successful vs. unsuccessful runs of the classic algorithms can be readily identified with automatic techniques, then we can train a self-supervisedly train a learning-based module to mimic and outperform the traditional method. This is a kind of knowledge distillation [19], from traditional models to deep ones.

Our optimization algorithm is essentially expectation maximization. After identifying the successful outcomes of a traditional object discovery method, we optimize the parameters of a learning-based method to make those successes more likely. We then use the combination of the traditional and learned methods to obtain new high-confidence estimates, and repeat. Our method outperforms not only the traditional methods, which only work under specific conditions, but also the deep methods, which currently only work in toy environments. We demonstrate success in a popular synthetic environment where recent deep models have been successful, and also on the real-world KITTI urban scenes benchmark, where the existing models fall flat.

Our main contribution is not in any particular component, but rather in their combination. We demonstrate that by exploiting the successful outcomes of simple handcrafted detectors and trackers, we can train a robust

108 learning-based method to detect and track objects in a target domain, without requiring any human annotations.
109
110

111 2. Related Work

112 **Object discovery** Many recent works have proposed deep
113 neural networks for object discovery. These models typi-
114 cally have an object-centric bottleneck, and are tasked with
115 a reconstruction objective. MONet [5], Slot attention [29],
116 IODINE [17] , SCALOR [22], AIR [11], and AlignNet [8]
117 fall under this category. These methods have been demon-
118 strated successful in a variety of simple domains, but have
119 not yet been tested on real-world videos. In this paper we
120 use the publicly available code for a subset of these mod-
121 els and evaluate whether they are still able to perform well
122 under the complexities of real-world imagery. Compared to
123 our method, these baselines do not make use of depth, but
124 in the supplementary file we show that concatenating depth
125 to their input stream provides very little impact on perfor-
126 mance.
127

128 **Ensemble methods** Using an ensemble is a well-known
129 way to improve the performance of a machine learning al-
130 gorithm. Assuming that each member of the ensemble is
131 prone to different types of errors, the combination of them
132 is likely to make fewer errors than any individual com-
133 ponent [9]. Ensembling is also the key idea behind knowledge
134 distillation, where knowledge gets transferred from cumber-
135 some models to simpler ones [19, 35]. The typical modern
136 setup is to make up the ensemble out of multiple copies of
137 a neural network, which are trained from different random
138 initializations and therefore have slightly different parame-
139 ters. In our case, the ensemble is more diverse: it is made
140 up of components which solve different tasks, but which can
141 still be checked against one another for consistency. For ex-
142 ample, we learn a 2D pixel labeller which operates on RGB
143 images, and a 3D object detector which operates on voxel-
144 ized pointclouds. When the 3D detections are projected
145 into the image, we expect them to land on “object” pixels.
146

147 **Never ending learning** We take inspiration from the
148 methods described in “never ending learning” [6]. We set
149 up a version of this in an expectation-maximization frame-
150 work, where in the expectation step we fire all modules and
151 look for agreement among them, and in the maximization
152 step we re-train the modules to improve this agreement.
153

154 **Structure-from-Motion/SLAM** Earlier works on struc-
155 ture from motion (SfM) methods [40, 7] set the ambitious
156 goal of extracting 3D scene pointclouds and camera 3D tra-
157 jectories from 2D pixel trajectories, exploiting the reduced
158 rank of the 2D trajectory matrix under rigid motions. Unfor-
159 tunately, these methods were often confined to very simple
160

161 videos, due to their difficulty handling camera motion de-
162 generacies, non rigid object motion and frequent occlusions
163 that caused 2D trajectories to be short in length. Simulta-
164 neous Localization And Mapping (SLAM) methods optimize
165 the camera poses in every frame as well as the 3D coor-
166 dinates of points in the scene online and in real time, and
167 assume a calibrated setup, i.e., that the focal length of the
168 camera is known [37, 24]. Our model exploits the rare suc-
169 cesses of these methods to learn about the static vs. moving
170 part of the scene.
171

172 **Moving object estimation** Most motion segmentation
173 methods cluster 2D optical flow vectors to segment moving
174 objects. While earlier approaches attempted motion seg-
175 mentation completely unsupervised by integrating motion
176 information over time through 2D pixel flow trajectories
177 [3, 32], recent works focus on learning to segment 2D ob-
178 jects in videos, supervised by annotated video benchmarks
179 [2, 20, 34, 26, 12, 25].
180

181 3. Track, Check, Repeat

182 Our method takes as input a video with RGB and depth
183 (either a depthmap or a pointcloud) and camera intrinsics,
184 and produces as output a 3D detector and 3D tracker for
185 salient objects in the video.
186

187 We treat this data as a test set, in the sense that we do not
188 use any annotations. In the current literature, most machine
189 learning methods have a training phase and a test phase,
190 where the model is “frozen” when the test phase arrives.
191 Our method instead attempts to optimize its parameters for
192 the test domain, using “free” supervision that it automati-
193 cally generates (without any human intervention).
194

195 The optimization operates in “rounds”. The first round
196 leverages optical flow and cycle-consistency constraints, to
197 discover moving objects in the videos. The most-confident
198 object proposals are upgraded into pseudolabels for train-
199 ing appearance-based object detectors. The second round
200 leverages optical flow and the new objectness detectors to
201 find more high-confidence proposals, which again lead to
202 additional training. In the final round we use the detectors
203 as trackers, and use these to generate a library of trajec-
204 tories, capturing a motion prior for objects in the domain.
205

206 A critical piece in each stage is the “check”, which de-
207 cides whether or not to promote an estimate into a pseudola-
208 bel for the next round. We now describe each piece of the
209 method, along with its corresponding check.
210

211 3.1. Optical flow

212 Optical flow indicates a 2D motion field that corresponds
213 to the pixels of a pair of images. We will use flow (in combina-
214 tion with other cues) as a signal for objectness, and also as a
215 submodule of egomotion estimation. As a starting point, we
216

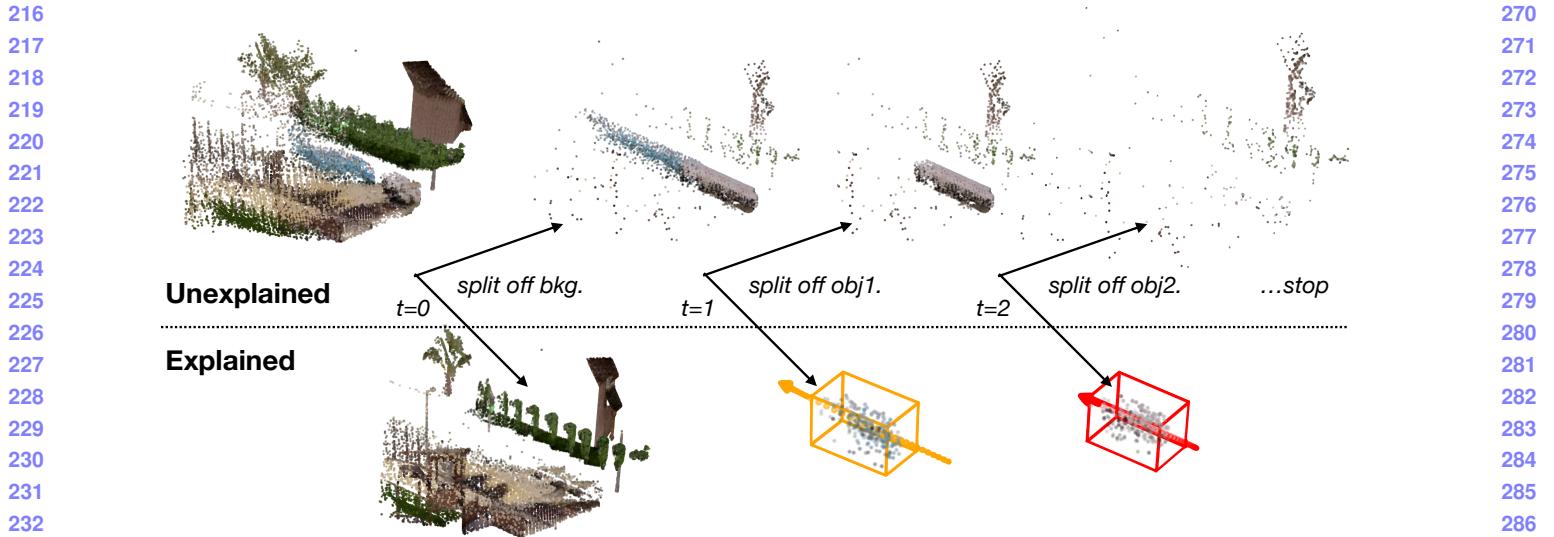


Figure 1. **Incrementally explaining an RGB-D video.** As input, our method consumes an RGB-D video, or alternatively, a set of RGB-D videos from a single test domain. This data is considered to be entirely “unexplained” at the start. This sequence is then “explained” in stages, which proceed from the most prominent moving components, which “pop out” due to motion saliency against the background, down to the less-prominent components which need to be discovered from appearance cues. Optimization proceeds in an expectation-maximization framework, where the parts of the scene that are already “explained” become pseudolabels for the next stage, which optimizes the parameters to explain more confidently with more modules, and generalize to the unexplained regions of the data.

use an off-the-shelf pre-trained convolutional optical flow network [39]. We note that optical flow can be learned entirely unsupervised [42].

We finetune this model with three self-supervision techniques. First, we use traditional edge-aware motion smoothness and brightness constancy [42]:

$$\mathcal{L}(\mathbf{w}; I_t, I_{t+1}) = \ell_{\text{photometric}}(\mathbf{w}; I_t, I_{t+1}) + \lambda \ell_{\text{smoothness}}(\mathbf{w}), \quad (1)$$

where $\{I_t, I_{t+1}\}$ are temporally consecutive images, $\mathbf{w} \in \mathbb{R}^{H \times W \times 2}$ is the optical flow, and λ is the parameter that weighs the relative importance of the smoothness. The photometric loss is computed by

$$\ell_{\text{photometric}}(\mathbf{w}; I_t, I_{t+1}) = \sum_{i,j} \rho_D(I_t - I_{t+1}(\mathbf{p} + \mathbf{w})), \quad (2)$$

where ρ_D is the robust generalized Charbonnier penalty function $\rho(x) = (x^2 + \epsilon^2)^\alpha$ to mitigate the effects of outliers. And the smoothness loss is computed by

$$\begin{aligned} \ell_{\text{smoothness}}(\mathbf{w}) = & \sum_j^H \sum_i^W [\rho_S(u_{i,j} - u_{i+1,j}) + \rho_S(u_{i,j} - u_{i,j+1}) \\ & + \rho_S(v_{i,j} - v_{i+1,j}) + \rho_S(v_{i,j} - v_{i,j+1})], \end{aligned} \quad (3)$$

where $u_{i,j}, v_{i,j}$ are horizontal and vertical components at (i, j) of \mathbf{w} , respectively, and $\rho_s(\cdot)$ is the spatial smoothness penalty function realized by the generalized Charbonnier function.

Second, we make a copy of the model and freeze it, to act as a “teacher” for another copy. From flows estimated by the teacher, we compute a mask indicating which flows are forward-backward consistent [41, 33, 38]. We first generate the forward flow $\mathbf{w}_{t \rightarrow t+1}$ and the backward flow $\mathbf{w}_{t+1 \rightarrow t}$. The reversed forward flow is computed by

$$\hat{\mathbf{w}}_{t \rightarrow t+1} = \mathbf{w}_{t+1 \rightarrow t}(\mathbf{p} + \mathbf{w}_{t \rightarrow t+1}(\mathbf{p})). \quad (4)$$

A pixel \mathbf{p} is considered to be forward-backward consistent if it satisfies

$$|\mathbf{w}_{t \rightarrow t+1} + \hat{\mathbf{w}}_{t \rightarrow t+1}|^2 < \alpha_1 \left(|\mathbf{w}_{t \rightarrow t+1}|^2 + |\hat{\mathbf{w}}_{t \rightarrow t+1}|^2 \right) + \alpha_2, \quad (5)$$

where we use $\alpha_1 = 0.05$ and $\alpha_2 = 0.5$ in our experiments.

This is a well-known “check” for optical flow – if flow is not cycle-consistent, it is not likely to be correct. We supervise the student model to mimic the teacher model at the pixels that are cycle-consistent. This helps ensure that the student model does not get worse than the teacher during optimization.

The third technique, inspired by the recent SelFlow method [28], is to apply random synthetic occlusions on the “student” model’s input. Synthetic augmentations on the student, while self-supervising via the teacher, forces the student model to learn a more robust flow estimator than the teacher.

324

3.2. Egomotion estimation

Egomotion is the rigid motion of the camera (i.e., transformation of poses) across a pair of frames. Estimating egomotion allows us to better estimate which pixels are moving due to the camera’s motion, and which are moving independently. Pixels moving independently are a strong cue for objectness.

We begin by “upgrading” the (dense) 2D flow fields into a sparse 3D pointcloud flow. To do this, we project the two pointclouds that correspond to the image pair, to obtain the nearest pixel coordinate of each 3D point. We then simply gather the flows whose startpoint and endpoint both have a corresponding 3D point, and use the delta of these points as the 3D pointcloud flow.

We then use RANSAC to estimate the rigid motion that explains the maximal number of 3D point flows. RANSAC is intended to be robust to outliers, but the answer returned is often catastrophically wrong.

The critical third step is to “check” the RANSAC output with a freely-available signal. The inlier count itself is such a signal, but this demands carefully tuning the threshold for inlier counting. Instead, we enforce cycle-consistency, similar to flow. We estimate the rigid motion twice: once using the forward flow, and once using the backward flow (which delivers an estimate of the inverse transform, or backward egomotion). We then measure the inconsistency of these results, by applying the forward and backward motion to the *same* pointcloud, and measuring the displacement:

$$XYZ'_0 = RT_{10}^{bw} RT_{01}^{fw}(XYZ_0) \quad (6)$$

$$err = \max_n(||XYZ'_0 - XYZ_0||), \quad (7)$$

where RT_{01}^{fw} denotes the rotation and translation computed from forward flow, which carries the pointcloud from timestep 0 to timestep 1, and RT_{10}^{bw} is the backward counterpart.

If the maximum displacement across the entire pointcloud is below a threshold (set to 0.25 meters), then treat the estimate as “correct”. In practice we find that this occurs about 10% of the time in the KITTI dataset.

On these successful runs, we use the egomotion to create another 3D flow field (in addition to the one produced by the flow module), and we subtract these to obtain the camera-independent motion field. Independently moving objects produce high-magnitude regions in the egomotion-stabilized motion field, which is an excellent cue for objectness.

In the first “round” of optimization, we proceed directly from this stage to pseudo-label generation. From the pseudo-labels, we then train the parameters of two object detectors, described next.

3.3. 2D objectness segmentation

This module takes an RGB image as input, and produces a binary map as output. The intent of the binary map is to estimate the likelihood that a pixel belongs to a moving object. This module transfers knowledge from the motion-based estimators into the domain of appearance, since it learns to mimic pseudolabels that were generated from motion alone. This is an important aspect of the overall model, since it allows us to identify objects from the “moving” type even when they are stationary.

We use a 50-layer ResNet [18] with a feature pyramid neck [27] as the architecture, and train the last layer with a logistic loss against sparse pseudo-ground-truth:

$$\mathcal{L}^{\text{seg}} = \sum \hat{m} \log(1 + \exp(-\hat{s} \cdot s)), \quad (8)$$

where m is a mask indicating where the supervision is valid. We experimented with and without ImageNet pretraining for the ResNet, and found that the pretrained version converges more quickly but is not very different in performance.

In training this module with sparse labels, it is critical to add heavy augmentations to the input, so that it does not simply memorize a mapping from the happenstance appearance to the sparse objectness labels. We use random color jittering, random translation and scaling, and random synthetic occlusions.

3.4. 3D object detection

This module takes as input a voxelized colorized pointcloud (computed from the RGB-D and intrinsics), and estimates oriented 3D bounding boxes of objects.

We use 3D U-Net-style convolutional encoder [36], and a CenterNet-style detection head [10]. The head produces a set of heatmaps, which encode objectness (in 1 channel), 3D size (in 3 channels), 3D subvoxel offset (in 3 channels), and orientation along the vertical axis (encoded as a categorical distribution over 16 channels). For implementation details we refer the reader to the supplementary file, and the original 2D CenterNet paper [10].

In training this module, we find that color augmentations have little effect, but randomized pointcloud orientation (when creating the voxelized input) is critical for learning an even distribution over possible object orientations. Additionally, we create partial occlusion augmentations, by randomly erasing an 20×20 area around the object center, while ensuring the target is not fully occluded.

3.5. Short-term tracking

To relocate a detected object over short time periods, we use the classic technique of cross correlation with a rigid template [30]. We do this using the features provided by the

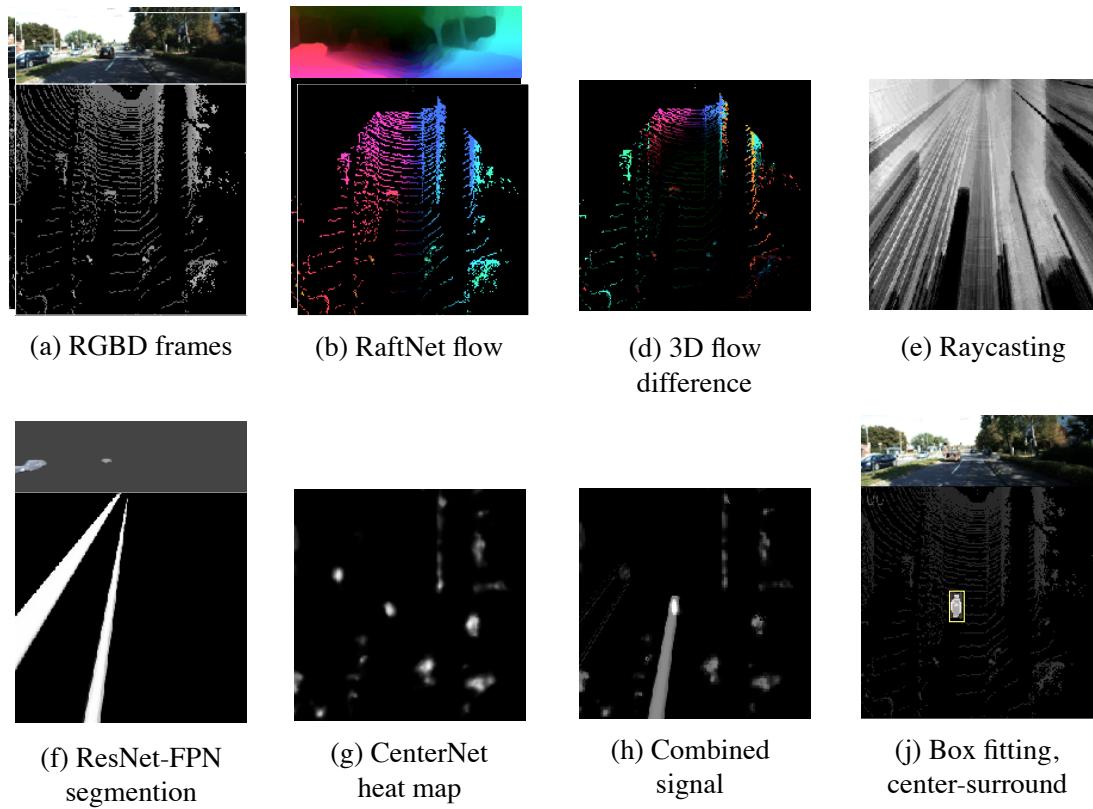
432
433
434
435
436
437
438
439
440
441
442
443
444
445446
447
448
449
450
451
452
453
454
455
456
457
458
459
460

Figure 2. Expectation (E) step of our unsupervised object discovery pipeline. Given an input video (a), we have multiple sources of evidence, including flow difference (c), visibility ray casting (e), 2D segmentation (f) and object-ness score from the previous round (g). Each of them can be error-prone, but ensembling them gives us high confidence object labels (j).

464

backbone of the object detector. We simply create a template by encoding a crop around the object, and then use this template for 3D cross correlation against features produced in nearby frames. We find that this is a surprisingly effective tracker despite not handling rotations, likely because the objects do not undergo large rotations under short timescales. To track for longer periods and across occlusions, we make use of motion priors represented in a library of previously-observed motions, described next.

474
475

3.6. Long-term tracking, with trajectory libraries

476
477
478
479
480

To track the objects, we need to associate the detection of the same objects through time. We use a trajectory library as a reference, to check the if the association is affordable, and to link detections across long temporal gaps (e.g., across occlusions).

481
482
483
484
485

To build the trajectory library, we first use nearest neighbour to link the detection from temporally consecutive frames. The trajectories are initialized by the detection in the first frame. Then we iterate through frames to extend the trajectories. We assume that the objects move with con-

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

stant velocity. Therefore, we can approximate the location of the object in the next frame $t + 1$. The detection that is close to the approximation and has similar size in frame $t + 1$ will be appended to the trajectory. The trajectories that last through the whole video will be included in the trajectory library since they are non-occluded and thus have high confidence to show the possible paths of objects.

During inference, we will first find the trajectories with the approach to build the trajectory library and mark the detections that have been associated to trajectories as explained. Notice that only those surviving through the whole video will be found. To also track the occluded objects, we use the trajectory library to link the unexplained detections. Given two unexplained detections D_{ts}, D_{te} from two different frames, where ts and te are the frame ids, we try to find the trajectories that have the same movement as the one from D_{ts} to D_{te} within $te - ts$ frames. Then we evaluate the possible trajectories in the library by computing the likelihood given the visual evidence. We generate the cost column by adding the centernet objectness heatmap, the unproject 2D segmentation heatmap and the invisible

540 map, where the centernet objectness heatmap and the un-project 2D segmentation heatmap provide the confidence
 541 that *visible* objects are on the trajectory and the invisible
 542 map shows where can have *invisible* objects that cannot be
 543 found by visual modules. The cost through the trajectory in
 544 the cost column is used as the likelihood. If there exist tra-
 545 jectories in the library that have high likelihood, we will say
 546 D_{ts} and D_{te} are likely to belong to the same trajectory, thus
 547 linking them as a new trajectory. Notice that the detections
 548 from other frames that are close to this trajectory will also
 549 be explained.
 550

552 3.7. Pseudo-label generation

553 Pseudo-label generation is what takes the model from
 554 one round of optimization to the next. The intent is to select
 555 the object proposals that are likely to be correct, and treat
 556 them as ground truth for training future modules.
 557

558 We take inspiration from never-ending learning architec-
 559 tures [31], which promote an estimate into a label only if (i)
 560 at least one module produces exceedingly-high confidence
 561 in the estimate, or (ii) multiple modules have reasonably-
 562 high confidence in the estimate.

563 The 2D and 3D modules directly produce objectness
 564 confidences, but the motion cues need to be converted into
 565 an objectness cue. Our strategy is inspired by classic liter-
 566 ature on motion saliency [21]: (1) compute the magnitude
 567 of the egomotion-stabilized 3D motion field, (2) threshold
 568 it at a value (to mark regions with motion larger than some
 569 speed), (3) find connected components in that binary map
 570 (to obtain discrete regions), and (4) evaluate the center-
 571 surround saliency of each region. Specifically, we com-
 572 pute histograms of the motion inside the region and in the
 573 surrounding shell, compute the chi-square distance between
 574 the distributions, and threshold on this value [1]:

$$575 \quad 576 \quad 577 \quad 578 \quad 579 \quad 580 \quad 581 \quad 582 \quad 583 \quad 584 \quad 585 \quad 586 \quad 587 \quad 588 \quad 589 \quad 590 \quad 591 \quad 592 \quad 593 \quad cs(\theta) = \chi^2(h(\text{cen}_\theta(\Delta XYZ)), h(\text{surr}_\theta(\Delta XYZ))), \quad (9)$$

where θ denotes the region being evaluated, cen_θ and surr_θ select points within and surrounding the region, h computes a histogram, and ΔXYZ denotes the 3D motion field.

When the trained objectness detectors are available (i.e., on rounds after the first), we convert the rigid motion field into a heatmap with $\exp(-\lambda ||\Delta XYZ||)$, and add this heatmap to the ones produced by the 2D and 3D objectness estimators, and then proceed with thresholding, connected components, and box fitting as normal. The only difference is that we use a threshold that demands multiple modules to agree: since each module produces confidences in $[0, 1]$, setting the threshold to any value above 2 effectively enforces this constraint.

4. Experiments

We test our model in the following datasets:

1. Synthetic RGB-D videos of tabletop scene rendered with CATER [16]. CATER is a built upon CLEVR [23], and it focuses on testing the model’s ability to do long term temporal reasoning. We modified the simulator so that it can generate pointclouds, but leave all other rendering parameters untouched. The max number of object is set to 10 to make the scene as complex as possible. Videos are captured by 6 virtual camera set static around the scene. 594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
2. Real RGB-D videos of urban scenes, from the KITTI dataset [15]. This data was collected with a sensor platform mounted on a moving vehicle, with a human driver navigating through a variety of areas in Germany. The data provides multiple RGB images; we use the “left” color camera. The dataset provides depth in the form of LiDAR sweeps synced to the RGB images. We use the “tracking” subset of KITTI, which includes 3D object labels, and approximate (but relatively inaccurate) egomotion provided by an inertial measurement unit. 604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

In all datasets, we test on 50-frame videos. We evaluate all models on their ability to discover objects in 3D. For our model and the traditional baseline, we additionally evaluate tracking performance.

4.1. Baselines

We evaluate the following unsupervised object discovery baselines. We also tested the baselines for unsupervised object tracking, but their performance is poor on both datasets, so we include the results in supplementary materials.

- **Object-Centric Learning with Slot Attention [29].** The Slot Attention model trains an autoencoder with image reconstruction loss. An iterative attention-based update mechanism is applied when constructing the slots, which act as a representational bottleneck. 628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
- **MONet [5].** MONet applies an recurrent attention mechanism that tries to explain the scene part by part. A VAE is also trained to reconstruct the image. Since there is no official code released, we re-implemented using PyTorch. 634
635
636
637
638
639
640
641
642
643
644
645
646
647
- **Discontinuities Tracing [4].** This method first extract dense point trajectories using optical flow. Video segmentation is then performed by detecting discontinuities of embedding density. We also tested a more naive approach **Spectral Clustering** [14] which the Discontinuities Tracing improved upon. 640
641
642
643
644
645
646
647

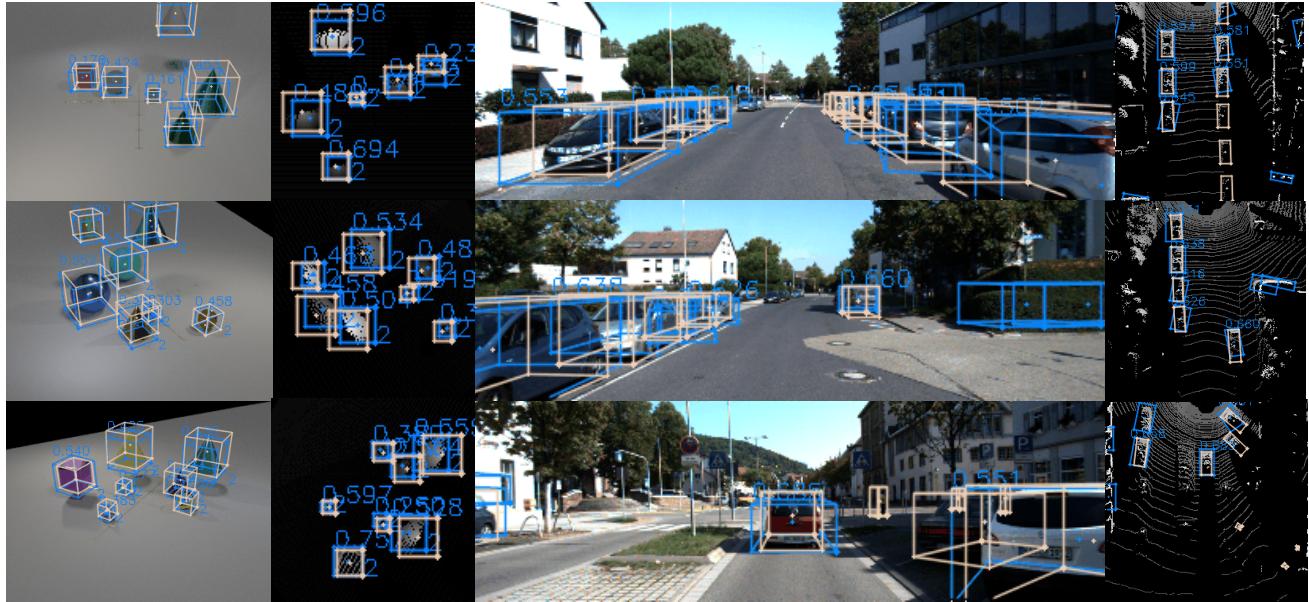


Figure 3. **3D object detections in CATER (left) and KITTI (right).** Ground-truth boxes are shown in beige and detection results are shown in blue. IoU scores are marked alongside each boxes. Results are shown in RGB and bird’s-eye view.

Table 1. Comparison of ours method to baselines performance on Object Discovery task. Results are reported as mAP at 7 different IoU thresholds. Ours method works best in all the matrices reported. 2D means perspective view and BEV means bird’s-eye view.

Method	Dataset	mAP@X						
		0.1	0.2	0.3	0.4	0.5	0.6	0.7
Slot Attention [29]	CATER (2D)	0.63	0.51	0.43	0.34	0.22	0.1	0.05
	KITTI (2D)	0.07	0.03	0.01	0	0	0	0
MONet [5]	CATER (2D)	0.23	0.14	0.12	0.10	0.07	0.03	0.01
	KITTI (2D)	0.03	0.01	0	0	0	0	0
Spectral Clustering [4]	CATER (2D)	0.18	0.08	0.04	0.03	0.01	0	0
	KITTI (2D)	0.01	0	0	0	0	0	0
Discontinuity Aware Clustering [14]	CATER (2D)	0.17	0.08	0.04	0.02	0.01	0.01	0
	KITTI (2D)	0.01	0	0	0	0	0	0
Ours (Round1)	CATER (2D)	0.98	0.97	0.97	0.94	0.86	0.7	0.36
	KITTI (2D)	0.53	0.39	0.18	0.06	0.03	0.01	0.01
	CATER (BEV)	0.97	0.92	0.75	0.57	0.34	0.06	0
	KITTI (BEV)	0.46	0.42	0.06	0	0	0	0
Ours (Round2)	CATER (2D)	0.98	0.97	0.96	0.94	0.88	0.69	0.33
	KITTI (2D)	0.43	0.40	0.39	0.33	0.30	0.22	0.10
	CATER (BEV)	0.97	0.95	0.84	0.66	0.46	0.08	0
	KITTI (BEV)	0.41	0.39	0.35	0.31	0.28	0.11	0.02

4.2. Quantitative Results

Object Discovery Our main evaluation is in Table 1, where we evaluate the object proposal accuracy in mean Average Precision (mAP) at different IoU threshold, on both CATER and KITTI dataset. The numbers are collected in Bird’s-eye view (BEV) and in RGB (2D). We find that our model produces the most accurate bounding boxes, in nearly all metrics. Adding another round of EM improves the numbers most of the time. Interestingly, the baselines methods that achieves state-of-the-art results in synthetic

datasets have near zero accuracy in KITTI, probably because their unsupervised learning objectives are not power enough to produce meaningful decomposition of complex scenes.

Object Tracking The object tracking accuracy (in IoU) is shown in Table 4. We provide the model with ground-truth bounding box at frame0. No ego-motion information are used during testing. Our methods can maintain a high IoU over a long time horizon. The IoU at frame19 is 0.34

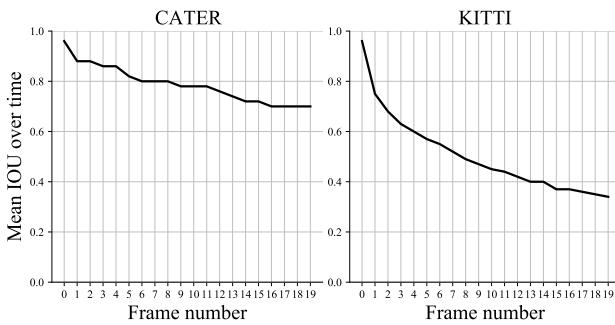


Figure 4. Object Tracking IOU in KITTI and CATER.

Table 2. Ablations of the trajectory library, in CATER.

Method	Accuracy	Precision
Ours, with nearest neighbour	0.53	0.94
... and trajectory library	0.64	0.91

in KITTI and 0.7 in CATER.

Ablation on the trajectory library The ablation study on the trajectory library is shown in Table 2. The accuracy is the ratio of the ground-truth trajectories that are found in predictions. The precision is the ratio of the predicted trajectories that successfully link the detections of the same objects. With the trajectory library, we improve the accuracy from 53% to 64%, while no big drop happens in precision measurements. It shows that our method can successfully track the objects through occlusion without blindly linking unexplained detections.

We present the remainder of the quantitative results in the supplementary, showing that the tracking performance of our model outperforms the baseline, and showing that ablating components of our model decreases performance.

4.3. Qualitative Results

For object discovery, We show object proposals of our model in CATER and KITTI in 3. Ground-truth boxes are shown in beige and proposed boxes are shown in blue. Their IoU are marked near the boxes. Results are shown on RGB image as well as bird’s-eye view. The boxes have high recall and high precision overall; it can detect small objects as well as separate the object that are spatially close to each other. In KITTI, there are some false positive results on bushes and trees because of the lack of pseudo-label supervision there.

We visualize object tracklets in KITTI in Figure 5, but we encourage the reader to inspect the supplementary video for clearer visualizations of the tracking performance.



Figure 5. 3D object tracking in KITTI. Ground-truth boxes are shown in beige and detection results are shown in blue. IoU scores are marked alongside each boxes. Each column shows timesteps from a video. Objects frequently undergo partial occlusions and truncation.

4.4. Limitations

The proposed method has three main limitations. Firstly, our work assumes access to RGB-D data with accurate depth, which is not affordable yet for many domains that our method has potential application on. Secondly, our model architecture requires a lot of GPU memory, due to its third spatial dimension. This severely limits either the resolution and speed. Finally, it is unclear for us how to mine for negatives. Right now we use a small region around each pseudo label as “negative”, but it leaves us open to false positives in far-away stuff like bushes and trees. Sparsifying our feature grid, or using points instead of voxels, are clear areas for future work.

5. Conclusion

We propose an unsupervised method that learns by segmenting motions of unlabelled RGB-D video streams. Learning and improvement is achieved in an expectation-maximization framework. In the expectation step we take the high confidence agreement among individual modules to create pseodo labels, and in the maximization step we re-train the modules. Standard data augmentation techniques are used to improve the generalization ability of the model. We showed that our method achieves state-of-the-art object discovery and tracking accuracy on two challenging datasets, and we further show that our method can address occlusions. Our approach opens new avenues for learning object models from videos in arbitrary environments, without requiring explicit object supervision. Extending our method to handle deformable and articulating motion is a useful avenue for future work.

864

References

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

- [1] Bogdan Alexe, Thomas Deselaers, and Vittorio Ferrari. Measuring the objectness of image windows. *IEEE transactions on pattern analysis and machine intelligence*, 34(11):2189–2202, 2012. [1](#), [6](#)
- [2] Goutam Bhat, Felix Järemo Lawin, Martin Danelljan, Andreas Robinson, Michael Felsberg, Luc Van Gool, and Radu Timofte. Learning what to learn for video object segmentation, 2020. [2](#)
- [3] Thomas Brox and Jitendra Malik. Object segmentation by long term analysis of point trajectories. In Kostas Daniilidis, Petros Maragos, and Nikos Paragios, editors, *ECCV*, pages 282–295, 2010. [2](#)
- [4] Thomas Brox and Jitendra Malik. Object segmentation by long term analysis of point trajectories. In *European conference on computer vision*, pages 282–295. Springer, 2010. [6](#), [7](#)
- [5] 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](#), [6](#), [7](#)
- [6] Andrew Carlson, Justin Betteridge, Bryan Kisiel, Burr Settles, Estevam R. Hruschka, and Tom M. Mitchell. Toward an architecture for never-ending language learning. In *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence*, AAAI’10, page 1306–1313. AAAI Press, 2010. [2](#)
- [7] J. Costeira and T. Kanade. A multi-body factorization method for motion analysis. *ICCV*, 1995. [2](#)
- [8] Antonia Creswell, Kyriacos Nikiforou, Oriol Vinyals, Andre Saraiva, Rishabh Kabra, Loic Matthey, Chris Burgess, Malcolm Reynolds, Richard Tanburn, Marta Garnelo, et al. Alignnet: Unsupervised entity alignment. *arXiv preprint arXiv:2007.08973*, 2020. [2](#)
- [9] Thomas G Dietterich. Ensemble methods in machine learning. In *International workshop on multiple classifier systems*, pages 1–15. Springer, 2000. [2](#)
- [10] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6569–6578, 2019. [4](#)
- [11] 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 D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29, pages 3225–3233. Curran Associates, Inc., 2016. [2](#)
- [12] Katerina Fragkiadaki, Pablo Arbelaez, Panna Felsen, and Jitendra Malik. Learning to segment moving objects in videos. In *CVPR*, June 2015. [2](#)
- [13] Katerina Fragkiadaki and Jianbo Shi. Exploiting motion and topology for segmenting and tracking under entanglement. In *CVPR*, 2011. [1](#)
- [14] Katerina Fragkiadaki, Geng Zhang, and Jianbo Shi. Video segmentation by tracing discontinuities in a trajectory em-

- bedding. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1846–1853. IEEE, 2012. [6](#), [7](#)
- [15] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013. [6](#)
- [16] Rohit Girdhar and Deva Ramanan. Cater: A diagnostic dataset for compositional actions and temporal reasoning. *arXiv preprint arXiv:1910.04744*, 2019. [6](#)
- [17] Klaus Greff, Raphaël Lopez Kaufman, Rishabh Kabra, Nick Watters, Chris Burgess, Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-object representation learning with iterative variational inference. *arXiv preprint arXiv:1903.00450*, 2019. [2](#)
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [4](#)
- [19] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. [1](#), [2](#)
- [20] Yuan-Ting Hu, Hong-Shuo Chen, Kexin Hui, Jia-Bin Huang, and Alexander G. Schwing. Sail-vos: Semantic amodal instance level video object segmentation - a synthetic dataset and baselines. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [2](#)
- [21] Laurent Itti and Christof Koch. A saliency-based search mechanism for overt and covert shifts of visual attention. *Vision research*, 40(10-12):1489–1506, 2000. [6](#)
- [22] Jindong Jiang*, Sepehr Janghorbani*, Gerard De Melo, and Sungjin Ahn. Scalor: Generative world models with scalable object representations. In *International Conference on Learning Representations*, 2020. [2](#)
- [23] 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, 2016. [6](#)
- [24] C. Kerl, J. Sturm, and D. Cremers. Dense visual SLAM for RGB-D cameras. In *IROS*, 2013. [2](#)
- [25] Anna Khoreva, Rodrigo Benenson, Eddy Ilg, Thomas Brox, and Bernt Schiele. Lucid data dreaming for object tracking. In *CVPR Workshops*, 2017. [2](#)
- [26] Yu Li, Zhuoran Shen, and Ying Shan. Fast video object segmentation using the global context module, 2020. [2](#)
- [27] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. [4](#)
- [28] Pengpeng Liu, Michael Lyu, Irwin King, and Jia Xu. Selflow: Self-supervised learning of optical flow. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4571–4580, 2019. [3](#)
- [29] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. *Advances in Neural Information Processing Systems*, 33, 2020. [2](#), [6](#), [7](#)

- 972 [30] Lain Matthews, Takahiro Ishikawa, and Simon Baker. The 1026
973 template update problem. *IEEE transactions on pattern 1027
974 analysis and machine intelligence*, 26(6):810–815, 2004. 4 1028
975 [31] Tom Mitchell, William Cohen, Estevam Rruschka, Partha 1029
976 Talukdar, Bishan Yang, Justin Betteridge, Andrew Carlson, 1030
977 Bhavana Dalvi, Matt Gardner, Bryan Kisiel, et al. Never- 1031
978 ending learning. *Communications of the ACM*, 61(5):103– 1032
979 115, 2018. 6 1033
980 [32] P. Ochs and T. Brox. Object segmentation in video: A 1034
981 hierarchical variational approach for turning point trajectories 1035
982 into dense regions. In *ICCV*, 2011. 2 1036
983 [33] Pan Pan, Fatih Porikli, and Dan Schonfeld. Recurrent 1037
984 tracking using multifold consistency. In *Proceedings of 1038
985 the Eleventh IEEE International Workshop on Performance 1039
986 Evaluation of Tracking and Surveillance*, 2009. 3 1040
987 [34] F. Perazzi, A. Khoreva, R. Benenson, B. Schiele, and A. 1041
988 Sorkine-Hornung. Learning video object segmentation from 1042
989 static images. In *2017 IEEE Conference on Computer Vision 1043
990 and Pattern Recognition (CVPR)*, pages 3491–3500, 2017. 2 1044
991 [35] Ilija Radosavovic, Piotr Dollár, Ross Girshick, Georgia 1045
992 Gkioxari, and Kaiming He. Data distillation: Towards omni- 1046
993 supervised learning. In *Proceedings of the IEEE Conference 1047
994 on Computer Vision and Pattern Recognition (CVPR)*, June 1048
995 2018. 2 1049
996 [36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U- 1050
997 net: Convolutional networks for biomedical image segmen- 1051
998 tation. In *International Conference on Medical image com- 1052
999 puting and computer-assisted intervention*, pages 234–241. 1053
1000 Springer, 2015. 4 1054
1001 [37] T. Schöps, J. Engel, and D. Cremers. Semi-dense visual 1055
1002 odometry for AR on a smartphone. In *ISMAR*, 2014. 2 1056
1003 [38] Ishwar K Sethi and Ramesh Jain. Finding trajectories of 1057
1004 feature points in a monocular image sequence. *IEEE Transac- 1058
1005 tions on pattern analysis and machine intelligence*, (1):56– 1059
1006 73, 1987. 3 1060
1007 [39] Zachary Teed and Jia Deng. Raft: Recurrent all- 1061
1008 pairs field transforms for optical flow. *arXiv preprint 1062
1009 arXiv:2003.12039*, 2020. 3 1063
1010 [40] Carlo Tomasi and Takeo Kanade. Shape and motion from 1064
1011 image streams under orthography: A factorization method. 1065
1012 *Int. J. Comput. Vision*, 9(2):137–154, Nov. 1992. 2 1066
1013 [41] Hao Wu, Aswin C Sankaranarayanan, and Rama Chellappa. 1067
1014 In situ evaluation of tracking algorithms using time reversed 1068
1015 chains. In *2007 IEEE Conference on Computer Vision and 1069
1016 Pattern Recognition*, pages 1–8. IEEE, 2007. 3 1070
1017 [42] Jason J. Yu, Adam W. Harley, and Konstantinos G. Derpanis. 1071
1018 Back to basics: Unsupervised learning of optical flow 1072
1019 via brightness constancy and motion smoothness. In *ECCV*, 1073
1020 2016. 3 1074
1021
1022
1023
1024
1025