

A Self-supervised Learning System for Object Detection in Videos Using Random Walks on Graphs

Juntao Tan, Changkyu Song and Abdeslam Boularias¹

Abstract—This paper presents a new self-supervised system for learning to detect novel and previously unseen categories of objects in images. The proposed system receives as input several unlabeled videos of scenes containing various objects. The frames of the videos are segmented into objects using depth information, and the segments are tracked along each video. The system then constructs a weighted graph that connects sequences based on the similarities between the objects that they contain. The similarity between two sequences of objects is measured by using generic visual features, after automatically re-arranging the frames in the two sequences to align the viewpoints of the objects. The graph is used to sample triplets of similar and dissimilar examples by performing random walks. The triplet examples are finally used to train a siamese neural network that projects the generic visual features into a low-dimensional manifold. Experiments on three public datasets, YCB-Video, CORe50 and RGBD-Object, show that the projected low-dimensional features improve the accuracy of clustering unknown objects into novel categories, and outperform several recent unsupervised clustering techniques.

I. INTRODUCTION

Robots are increasingly deployed in challenging environments that contain unknown objects. Examples of such environments include households, warehouses and workshops, where robots are tasked with picking specific items from dense piles of a large variety of objects [1], [2]. Current robotic systems solve this problem by using a convolutional neural network (CNN) for detecting objects in images. CNNs are typically trained by using a large number of manually labeled images, which is a tedious process [3]–[5]. In this work, we propose a new self-supervised system that allows robots to learn novel categories of encountered objects on their own.

The proposed system receives several videos of piles of various unknown objects. Consecutive frames in the videos are obtained by randomly moving the camera or the objects to expose different viewpoints. The videos can be recorded at different times or in different locations. Therefore, the relation between the objects in the different videos is completely unknown. The frames may also contain various types of objects that belong to the same category, such as different types of coffee mugs for example. The goal of the robot is to autonomously: *i*) segment each frame into objects, *ii*) cluster the objects from all the frames and sequences into categories, *iii*) assign a numerical label to each discovered category, and *iv*) train a CNN using the automatically labeled data to

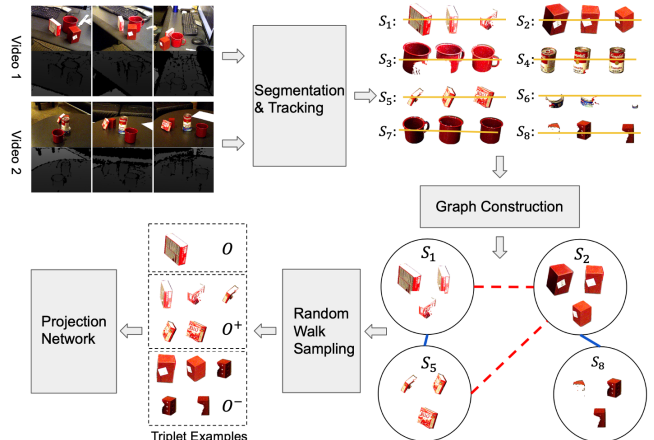


Fig. 1. Overview of the proposed system

recognize the newly discovered categories in future images. We focus in this work on steps *ii* and *iv*, and utilize for step *i* the technique presented in [6] for unsupervised segmentation using depth information. The last step, *iv*, depends only on the accuracy of the labels generated by the proposed system. Therefore, we focus in this paper on assessing the accuracy of the proposed unsupervised object clustering process.

The proposed approach builds on top of recent *self-supervised* techniques that utilize siamese neural networks with triplet loss functions to learn visual representations [7]–[14]. The networks are trained to project images, or high-dimensional features extracted from pre-trained networks such as *ResNet*, into low-dimensional feature vectors that can then be used for various tasks, such as clustering. The triplet ranking loss is designed such that the projected features of similar objects, i.e. objects that belong to the same category, are closer to each other than to the projected features of dissimilar objects, i.e. objects from other categories. The key challenge is finding a large number of examples of similar and dissimilar objects without supervision. While most existing techniques were designed for static images, some methods use 2D tracking in videos to automatically generate a large number of examples of the same object from different viewpoints [7], [8]. Examples of dissimilar objects are obtained by randomly selecting frames from different videos. There are two main issues in this approach. First, tracking can only provide examples of the same object. Other objects that belong to the same category typically appear in different videos. Second, randomly selected frames can possibly contain objects that belong to the same category and that cannot thus be used as examples of dissimilar objects.

¹The authors are with the Department of Computer Science of Rutgers University, Piscataway, New Jersey 08854, USA. {jt867, cs1080, ab1544}@cs.rutgers.edu. This work is supported by NSF awards IIS-1734492 and IIS-1846043.

The main contribution of our work is a new solution to these two issues that consists in constructing a graph where the nodes are sequences of object masks in consecutive frames, and the weights of the edges are the degrees of similarities between the sequences. The degree of similarity between two sequences is computed by searching for the best alignment of the different viewpoints of the objects in the two sequences. The degrees of similarities are interpreted as transition probabilities. A random walk on the graph is then used to generate examples of similar and dissimilar objects from different videos. Extensive evaluations on three publicly available video datasets clearly show that the learned features can effectively be used to cluster without supervision novel objects according to their unknown semantic labels. The proposed method also outperforms several clustering and self-supervised representation learning techniques.

II. RELATED WORK

Self-supervised learning of visual representations is becoming increasingly popular due to the colossal manual labeling efforts required by traditional deep learning techniques [15]–[17]. For example, it has been shown in [7] that efficient features can be learned from unsupervised auxiliary tasks, such as *context prediction*. In a closely related work [9], a triplet loss function is used to learn visual representations from videos. Unlike in the proposed method, negative examples in [9] are selected randomly while assuming that other videos contain only categories of objects other than that of the anchor patch. Moreover, the objective of [9] is learning feature descriptors that are then used for supervised classification tasks with labeled examples, which is different from our objective. A triplet-siamese network was also used for unsupervised visual representation learning in [8], where the triplet examples are obtained from simple transitive relations in a similarity graph. As in our proposed approach, intra-instance variations are obtained by tracking an object in a video sequence. While the approach proposed in [8] relies on the prior work [7] to find inter-instance invariances, our approach uses a more appropriate measure of similarity that is based on clustering objects into a large number of viewpoints in each video sequence, and then solving an assignment problem that matches viewpoints taken from different sequences. Moreover, only one-step transitive relations are considered in [8], while our approach utilizes a long-horizon random walk in the graph to sample positive examples by interpreting distances as inverse transition probabilities. Finally, negative examples are sampled randomly in [8], whereas they are sampled in our method from the complementary random walk distribution.

Other works on self-supervised learning from images construct image representations that are semantically meaningful via pretext tasks that do not require semantic labeling [11], [18]–[23]. For example, the Pretext-Invariant Representation Learning (PIRL) [18] approach learns invariant representations by using pretext task that involves solving jigsaw puzzles. This approach was designed for static images.

Our approach achieves similar objectives for videos. Earlier works on unsupervised learning of invariant features from videos [24]–[26] were proposed prior to [9], but they were also limited to tracking objects within a single sequence. Various techniques for clustering image features have been used in the past for detecting object categories without labels [27]–[44]. These techniques however rely on pre-trained features without fine-tuning them to improve the task of categorizing novel objects in a given small set of images encountered by a robot, which is our main objective.

III. PROPOSED APPROACH

A. Problem Setup

We consider the following problem. There is a set $\mathcal{V} = \{V_1, V_2, \dots, V_m\}$ of m video sequences. Each video sequence has a maximum of n frames. Each one of the frames is denoted by $f_t^{(i)}$ where t is the time of the frame in sequence i . Therefore, $V_i = (f_1^{(i)}, f_2^{(i)}, \dots, f_n^{(i)})$. Suppose now that there are K semantic classes $\{c_1, c_2, \dots, c_K\}$ of objects that appear in different or same frames. Examples of semantic classes are mugs, bowls, scissors, *etc.* Suppose that there are N individual instances of objects that appear in these frames. The number of instances is larger than or equal to the number of classes, i.e. $N \geq K$. The problem consists in segmenting each frame into individual objects, and then clustering all the segments from the different frames and sequences into K clusters, such that each cluster contains only objects that have the same semantic label. The challenge for the robot is to perform this task without any external supervision, so that the robot can self-label images of new types of objects and use the automatically labeled images to train object detectors, in a lifelong learning process that does not require human assistance. This challenge is exacerbated by the fact that objects belonging to the same class typically have different shapes and colors (*inter-instance variations*), and the same object appears in different frames with different viewpoints, illuminations and occlusions (*intra-instance variations*).

B. Overview

Figure 1 depicts an overview of our proposed system. It consists in *i*) segmenting RGB-D frames into individual objects, *ii*) tracking the objects along each sequence of frames, *iii*) clustering different viewpoints of each tracked object into a small number to reduce the number of frames, *iv*) measuring similarities between different sequences by solving an optimal assignment problem between viewpoints, *v*) performing random walks on the graph of similarities to generate similar and dissimilar examples, and finally *vi*) using the self-generated examples to learn a projection of visual features into a low-dimensional manifold with a soft triplet loss function. Features of objects in the low-dimensional manifold are clustered into K clusters by using k -means. These steps are explained in the following.

C. Segmentation and Tracking of Individual Instances

1) Segmentation: We follow the unsupervised segmentation approach that we have previously proposed in [6]. This

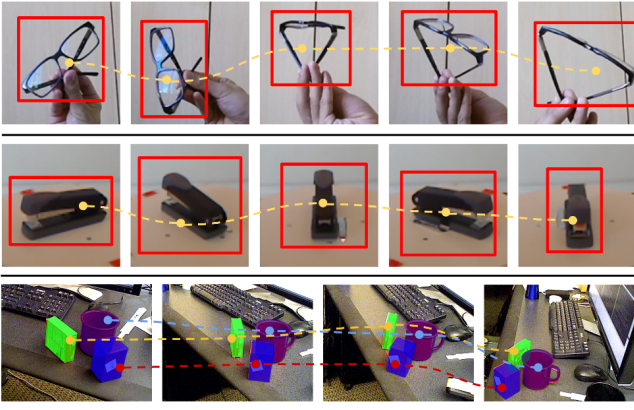


Fig. 2. Examples of three sequences of object masks obtained from segmentation and tracking in the datasets CORE50 [45], [46] (top row) RGBD-Object [47] (middle row), and YCB-Video [48] (bottom row).

approach transforms the point cloud of a frame $f_t^{(i)}$ into a graph of nearest neighbors, and utilizes the spectral clustering technique to merge supervoxels in the point cloud into segments. This method makes one assumption: the surface of an object is overall approximately convex. Consequently, certain objects may be over-segmented into convex parts. We solve this problem by using the color watershedding technique to merge the over-segmented objects in a post-processing step. The segmentation module returns a set $\mathcal{O}_t^{(i)}$ of object masks for each frame $f_t^{(i)}$ in each video V_i .

2) *Tracking*: In each video V_i , object masks $O_{t,j}^{(i)} \in \mathcal{O}_t^{(i)}$ at different times t that correspond to the same instance are linked together by tracking the motion of the masks over time. We start by creating a 2D bounding box around each object mask $O_{t,j}^{(i)} \in \mathcal{O}_t^{(i)}$ in each frame $f_t^{(i)}$ of video V_i . We assume that the frame rate of the videos is sufficiently high so that the 2D bounding boxes of the same object in any two consecutive frames $f_t^{(i)}$ and $f_{t+1}^{(i)}$ overlap. We denote the Intersection over Union (IoU) of the bounding boxes of masks $O_{t,j}^{(i)}$ and $O_{t+1,k}^{(i)}$ by $IoU(O_{t,j}^{(i)}, O_{t+1,k}^{(i)})$. The following linear matching problem is then solved for every video sequence V_i and every frame $f_t^{(i)}$,

$$\begin{aligned} & \text{Minimize } \sum_j \sum_k M(j, k) IoU(O_{t,j}^{(i)}, O_{t+1,k}^{(i)}) \\ & \text{s.t. } \forall j : \sum_k M(j, k) = 1, \forall k : \sum_j M(j, k) = 1, \\ & \quad \forall j, k : 0 \leq M(j, k) \leq 1. \end{aligned}$$

When $M(j, k) = 1$, object masks $O_{t,j}^{(i)}$ and $O_{t+1,k}^{(i)}$ are considered as masks of the same object, taken at consecutive times t and $t+1$ within video sequence V_i .

By arranging consecutive masks of the same object into a sequence, the final result of the segmentation and tracking process is a set $\mathcal{S} = \{S_1, S_2, \dots, S_l\}$ of mask sequences, where $S_i = (O_1^{(i)}, O_2^{(i)}, \dots, O_h^{(i)})$ is a sequence of masks $O_t^{(i)}$ that correspond to the same instance of object within a video. Note that each frame in a video can contain multiple objects. Thus, the same video sequence can yield several mask sequences, one for each tracked object. Examples of obtained mask sequences S_i are illustrated in Figure 2.

In the following, we show how to measure the likelihood

that two sequences S_i and S_j of masks correspond to the same class of objects. We then utilize this similarity measure to construct a graph and exploit the transitive relations in this graph, through random walks, to construct training examples for learning the visual features that will be used for categorizing individual masks $O_t^{(i)}$ into K clusters.

D. Similarity Graph Construction

A graph of inter-instance invariances is defined by the set of vertices \mathcal{S} , wherein each vertex $S_i \in \mathcal{S}$ is a sequence of a tracked individual object in a video sequence. The weight $W(S_i, S_j)$ of an edge (S_i, S_j) measures how likely are objects tracked in S_i and S_j to belong to the same class. This problem is highly challenging since we do not assume to have access to the list of semantic classes nor to any labeled data. Weights $W(S_i, S_j)$ are defined as follows,

$$W(S_i, S_j) = \max\left(\lambda W^+(S_i, S_j) - W^-(S_i, S_j), 0\right),$$

wherein $W^+(S_i, S_j)$ measures the similarity between sequence S_i and sequence S_j , $W^-(S_i, S_j)$ measures the dissimilarity between the two sequences, and λ is a constant hyper-parameter factor. The remainder of this subsection explains how $W^+(S_i, S_j)$ and $W^-(S_i, S_j)$ are computed.

1) *Computing W^+* : We start by extracting generic visual feature vectors $\Phi(O_t^{(i)})$ for every object mask $O_t^{(i)}$ in every frame-time t and every mask sequence S_i . Any standard feature extractor, such as HOG, SIFT or ResNet pre-trained offline on different types of objects and images, can be used for this purpose. We then cluster all the feature vectors from all the frames and object masks into a large number of *global clusters* by using the k -means algorithm. The number of global clusters is so large (e.g., $k = 500$ clusters in our experiments) that each cluster contains only a few objects. Objects belonging to the same global cluster are thus highly likely to belong to the same semantic class, and they are often the same instance seen from different viewpoints. Similarity weight $W^+(S_i, S_j)$ is simply the number of pairs of feature vectors $\Phi(O_t^{(i)})$ and $\Phi(O_{t'}^{(j)})$ that appear in the same cluster.

2) *Computing W^-* : $W^-(S_i, S_j)$ measures the distance between sequences S_i and S_j . To obtain this distance, one cannot simply add together the distances between features vectors $\Phi(O_t^{(i)})$ and $\Phi(O_{t'}^{(j)})$ of object masks $O_t^{(i)}$ and $O_{t'}^{(j)}$ at the same time-frames t in the two sequences, because the viewpoints in the two sequences are arbitrary and unaligned. Thus, a linear optimal assignment problem needs to be solved here in order to align the two sequences as well as possible by re-arranging their frames. Note that the two sequences do not necessarily correspond to the same class of object. For example, sequence S_i could be tracking a mug, while S_j is tracking scissors. In that case, re-arranging the frames to align viewpoints in the two sequences is futile. However, the resulting distance would be higher than the distance between S_i and another sequence S_k that tracks the same category of object, such as a different mug, as illustrated in Figure 3. Distance $W^-(S_i, S_j)$ is defined as,

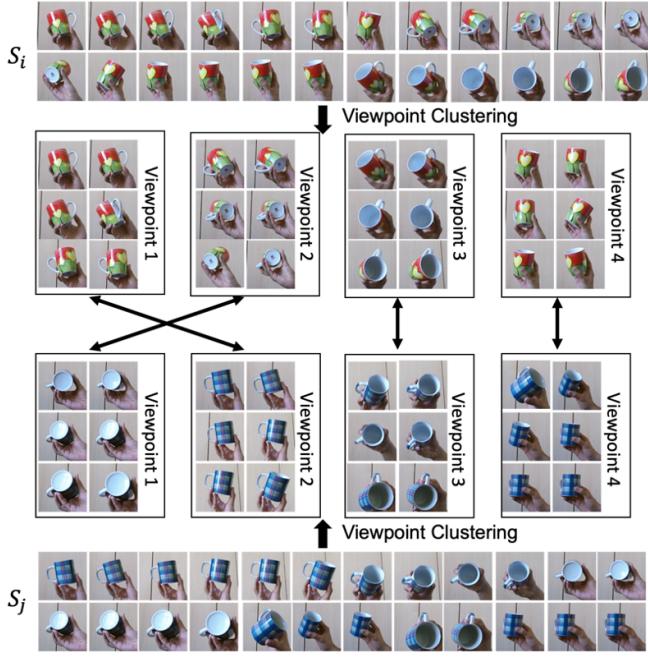


Fig. 3. Viewpoint matching

$$W^-(S_i, S_j) = \min_M \sum_t \sum_{t'} M(t, t') \|\Phi(O_t^{(i)}) - \Phi(O_{t'}^{(j)})\|_2$$

$$\text{s.t. } \forall t : \sum_{t'} M(t, t') = 1, \forall t' : \sum_t M(t, t') = 1,$$

$$\forall t, t' : 0 \leq M(t, t') \leq 1.$$

When $M(t, t') = 1$, object masks $O_t^{(i)}$ and $O_{t'}^{(j)}$ at frame-times t and t' in sequences S_i and S_j respectively are considered to belong to the same class of objects and to have the same viewpoint.

Since the frame rate of the videos is typically high, the viewpoint matching process is computationally expensive. To address this issue, we do not compare individual feature vectors directly. Feature vectors are first clustered into a small number of viewpoint regions. Feature vectors Φ in the objective function of the optimization problem above are substituted by centroids of their clusters, as illustrated in Fig. 3. This optional step not only reduces the computational cost of the viewpoint matching process, but it also ensures a balance between different viewpoints. For example, the camera may focus on a certain angle of the scene for a long period of time before moving to a different angle. Therefore, features of objects taken from the first angle will be over-represented in the sequence. Clustering overcomes this issue and ensures that different angles contribute equally to the objective function.

After computing W^+ and W^- , we can compute edge weights W . Figure 4 shows a concrete example of the resulting similarity graph from our experiments. In the constructed graph, every node is a tracking sequence $S_i = (O_1^{(i)}, O_2^{(i)}, \dots, O_h^{(i)})$. Two nodes S_i and S_j are connected only when their weight $w(S_i, S_j)$ is strictly positive. Thus, increasing the value of λ results in denser graphs.

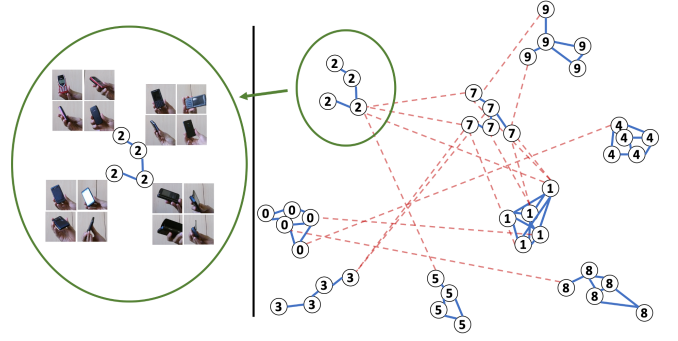


Fig. 4. Similarity Graph. This graph was generated by our system for the dataset CORE50 [45], [46] without using labeled data or supervision. Each node is a sequence of tracked object masks. The number inside each node indicates the ground-truth class of the object. The lengths of the edges are the inverse of their weights. Dashed edges indicate false similarities, notice how they are longer than the positive similarities (blue lines).

E. Sampling Triplet Examples

The constructed similarity graph loosely indicates which objects belong to the same class. But the graph typically contains several inaccuracies due to large inter-instance and intra-instance variations, and the fact that the relations in the graph are primarily extracted by using generic visual features along with temporal information in the viewpoint matching process. We demonstrated in our experiments that a final step is necessary for improving the accuracy of the main task of clustering objects into semantic classes. This step consists in creating a pool of triplets of examples $\langle O, O^+, O^- \rangle$ where O and O^+ are examples of hypothetically similar objects (same class), and O and O^- are examples of hypothetically dissimilar objects. The examples are then used to train a siamese-triplet network to project generic feature vectors Φ into a low-dimensional manifold, as illustrated in Fig. 5.

We first explain in the following how triplet examples are mined from the graph, then we show how the siamese-triplet network is trained using the generated examples. Each row in the weight matrix W of the graph is normalized by dividing each entry $W(S_i, S_j)$ by $\sum_k W(S_i, S_k)$. The resulting matrix, denoted by T , is a stochastic transition matrix. In other terms, $T(S_i, S_j)$ is the probability of selecting an object in sequence S_j as a positive example of objects in sequence S_i . As we can see from Figure 4, vertices corresponding to sequences of the same class of objects tend to be clustered together. To take full advantage of this information, we consider transitive relations between the vertices. We thus start at a node S_i and perform a random walk on the graph to sample a similar example S_j . For efficiency, we compute a probability distribution of the visitation frequencies, and use it to select similar and dissimilar examples. This distribution is given by the matrix T^H , which is computed recursively as $T^1 = T$ and $T^{H+1} = TT^H$. Examples that are similar to S_i are sampled from $T^H(S_i, \cdot)$, dissimilar examples are sampled from the complementary normalized distribution $\left(1 - T^H(S_i, \cdot)\right) \left(\left(1 - T^H(S_i, \cdot)\right) \mathbf{1}^T\right)^{-1}$ where $\mathbf{1}$ is a row vector with all elements equal to 1.

The result of the sampling process is a set of several

triplets $\langle S, S^+, S^- \rangle$, wherein (S, S^+) are similar sequences and (S, S^-) are dissimilar sequences. We sample from each trajectory several frames, and the result is a set of triplets $\langle O, O^+, O^- \rangle$, wherein (O, O^+) are similar objects and (O, O^-) are dissimilar objects. These examples are used to train a siamese-triplet network that projects features $\Phi(O)$ of objects into low-dimensional features $\Gamma(\Phi(O))$. We propose the following soft loss to train the network,

$$L_{\Gamma}(O, O^+, O^-) = \max \left(\|\Gamma(\Phi(O)) - \Gamma(\Phi(O^+))\|_2 - \|\Gamma(\Phi(O)) - \Gamma(\Phi(O^-))\|_2 + \alpha \text{conf}(O, O^+, O^-), 0 \right),$$

where α is a hyper-parameter and conf is defined as $\text{conf}(O, O^+, O^-) = \min \left(\frac{W(S, S^+) - \min_{S'} W(S, S')}{\max_{S'} W(S, S') - \min_{S'} W(S, S')}, 1 - \frac{W(S, S^-) - \min_{S'} W(S, S')}{\max_{S'} W(S, S') - \min_{S'} W(S, S')} \right)$, where S, S^+, S^- are the sequences from which O, O^+, O^- are respectively taken. $\text{conf}(O, O^+, O^-)$ is a number between 0 and 1 that indicates the confidence in O and O^+ belonging to the same class, and O and O^- belonging to different classes.

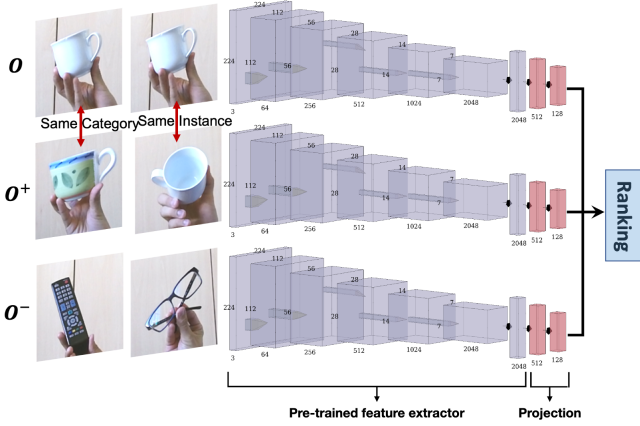


Fig. 5. Projection Network

IV. EXPERIMENTS

A. Datasets

The proposed method is evaluated on three public datasets that are designed for robotic tasks: RGBD-Object [47], CORE50 [45], [46] and YCB-Video [48]. RGBD-Object dataset contains 300 object instances classified into 51 categories. CORE50 has 50 object instances from 10 categories. YCB-Video contains 21 objects from 21 categories.

B. Setting

After training the projection layers of the network (Fig. 5), we simply apply K -means on the projected features $\Gamma(\Phi(O))$, where K is the number of classes. For the YCB-Video dataset, we set $K = 20$ because the ‘large clamp’ and ‘extra large clamp’ are the same object with slightly different sizes. Objects in the YCB-Video dataset are highly occluded, which makes the segmentation more challenging than in the two other datasets. Therefore, we also test a variant of our system on the YCB-Video dataset where

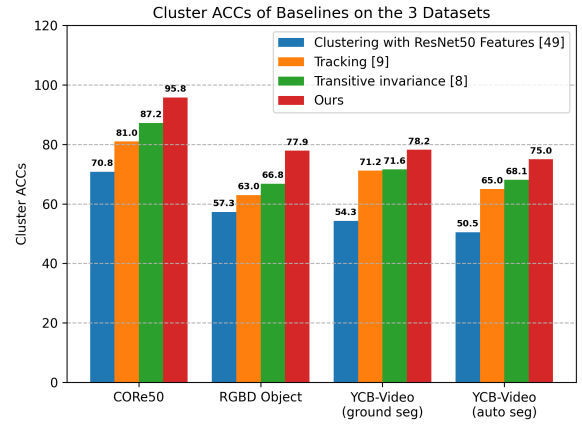


Fig. 6. Average Cluster Accuracies (ACC) of the compared techniques

the automatic segmentation of [6] is replaced with ground-truth segmentation of point clouds into individual objects, in order to show the potential of the proposed system if the RGB-D segmentation is improved, since our approach is independent of the segmentation method. The number of viewpoint clusters (Fig. 3) is set to 5 in all our experiments. For CORE50, RGBD-Object and YCB-Video datasets respectively, we use $\lambda = 0.1$, $\lambda = 1$ and $\lambda = 0.1$ to construct the similarity graph and horizons $H = 3$, $H = 5$, and $H = 3$ to sample triplets with random walks. We use *ResNet50* for extracting the high-dimensional generic features $\Phi(O)$. The projection layers that map $\Phi(O)$ into low-dimensional features $\Gamma(\Phi(O))$ are two fully connected layers with 512 and 128 output dimensions. ReLu activation function is used after the first projection layer. We fix the feature extractor parameters and only train the projection layers. The learning rate is set to 0.01, the margin value α is set to 10, and we use the stochastic gradient descent (SGD) optimizer.

C. Compared methods

We compare our method against the following alternative techniques: **1) ResNet+k-means** [49]: We use k -means directly on the features returned by ResNet50 pre-trained on *ImageNet*, where k is set to the number of categories in each dataset. **2) Tracking** [9]: Projection layers on top of ResNet50 are trained with positive examples sampled only from the same video sequence, and negative examples sampled randomly from other videos. **3) Transitive Invariance** [8]: An alternative graph-based approach for mining negative and positive examples for training the same projection layers on top of pre-trained ResNet50. We also compare against the following deep clustering techniques: **4) Deep Embedded Clustering (DEC)** [33], and **5) Deep Clustering** for unsupervised learning of visual features [42].

D. Results

We use three metrics to evaluate the final clustering results: 1) Average Cluster Accuracy (ACC), 2) Adjusted Rand Index (ARI) score and 3) Normalized Mutual Information (NMI) score. Results reported in Table I and illustrated in Figure 6 show that the proposed system significantly outperforms the other baselines. Figure 7 illustrates examples of the clusters

	COrE50			RGBD-Object			YCB-Video (ground-seg)			YCB-Video (auto-seg)		
	ACC	ARI	NMI	ACC	ARI	NMI	ACC	ARI	NMI	ACC	ARI	NMI
ResNet+K-Means [49]	70.8%	0.695	0.765	57.3%	0.456	0.726	54.3%	0.499	0.698	50.5%	0.345	0.586
DEC [33]	79.2%	0.776	0.863	56.5%	0.449	0.754	53.7%	0.559	0.764	45.6%	0.359	0.591
Deep Cluster [42]	51.1%	0.348	0.493	41.4%	0.309	0.624	45.8%	0.324	0.539	41.6%	0.254	0.464
Tracking [9]	81.0%	0.830	0.898	63.0%	0.551	0.808	71.2%	0.630	0.784	65.0%	0.517	0.692
Transitive Invariant [8]	87.2%	0.822	0.885	66.8%	0.572	0.817	71.6%	0.611	0.775	68.1%	0.579	0.707
Ours (Binary Graph)	93.7%	0.914	0.923	71.1%	0.650	0.844	72.2%	0.631	0.793	70.5%	0.589	0.725
Ours (Weighted Graph)	95.8%	0.927	0.953	77.9%	0.742	0.889	78.2%	0.682	0.810	75.0%	0.624	0.741

TABLE I



Fig. 7. Examples of clusters returned by the proposed self-supervised system (one per row and per dataset). The system succeeded in discovering the different categories of objects in the videos, and in assigning each object into a cluster that contains mostly only objects of the same category. Some objects (shown in red boxes) are misclassified. The discovered categories are named numerically by the system. For example, ‘label 4’ in YCB-Video dataset refers to ‘power drills’.

returned by our approach. Interestingly, most of the clusters contain only objects that belong to the same semantic class when the number of clusters is set to the number of classes. Similar levels of accuracy are observed even when the objects are highly occluded, such as in the YCB-video dataset.

E. Ablation Studies

We performed ablation studies to evaluate the impact of different aspects of our approach. In the first study, we tested our method against a variant where the weighted graph is replaced by a **binary graph** that preserves only the structure of the original graph but not the weights. Table I shows that the weights, computed by measuring similarities between sequences, play a major role in the performance of the system. In the second study, we trained the projection network using the standard triplet loss, i.e. without using the confidence function and with the margin term set to a constant $\alpha = 10$. The ACC results are 94.1%, 73.1% and 72.0% for COrE50, RGBD-Object and YCB. These results are below the ones obtained by our proposed soft triplet loss using the confidence function (last row in Table I).

Finally, we tested various approaches for measuring similarities between two sequences. The first one simply averages the feature vectors of all frames in a sequence and returns its distance from the mean feature vector of the second sequence. The second method clusters frames into viewpoints, like in our method, but does not align viewpoints with the compared sequence. Instead, it returns the average distance of the top ten nearest pairs of viewpoints (one from each

Matching Methods	COrE50	RGBD-Object	YCB-Video
Mean Feature Distance	0.75	0.61	0.52
Top Ten Nearest Neighbors	0.70	0.55	0.53
Cut Sequence Matching	0.79	0.61	0.56
Viewpoint Matching (ours)	0.82	0.65	0.56

TABLE II

sequence). The last method cuts each video evenly into ten parts, and returns the average distance between the means of the parts from the two sequences. To evaluate these methods, we match two sequences if their distance is smaller than a threshold and evaluate the matching quality with the f_β score defined as $f_\beta = (1 + \beta^2) \cdot \frac{\text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$. We set $\beta = 0.5$. Results reported in Table II clearly show that the proposed matching technique achieves the best results, especially in the COrE50 dataset where objects are rotated by 2π in each video. Results, videos, code, and data are available at <https://github.com/chrisjtan/RWS>.

V. CONCLUSION

Self-supervised object detection and recognition is an important skill that robots need to acquire on the road towards sustainable full autonomy. We have shown in this paper that such skills can be acquired by using robust tracking and matching techniques that take advantage of rich information contained in videos, along with the transitive nature of object similarities. In a future work, we plan to utilize the numerical labels that are automatically generated by our system to train an FCN for semantic segmentation in order to quickly detect the same types of objects in future images, without the need to run our entire pipeline.

REFERENCES

- [1] C. Mitash, R. Shome, B. Wen, A. Boularias, and K. E. Bekris, "Task-driven perception and manipulation for constrained placement of unknown objects," *IEEE Robotics Autom. Lett.*, vol. 5, no. 4, pp. 5605–5612, 2020. [Online]. Available: <https://doi.org/10.1109/LRA.2020.3006816>
- [2] R. Shome, W. N. Tang, C. Song, C. Mitash, H. Kourtev, J. Yu, A. Boularias, and K. E. Bekris, "Towards robust product packing with a minimalistic end-effector," in *International Conference on Robotics and Automation, ICRA 2019, Montreal, QC, Canada, May 20-24, 2019*. IEEE, 2019, pp. 9007–9013. [Online]. Available: <https://doi.org/10.1109/ICRA.2019.8793966>
- [3] A. Zeng, K.-T. Yu, S. Song, D. Suo, E. Walker Jr, A. Rodriguez, and J. Xiao, "Multi-view self-supervised deep learning for 6d pose estimation in the amazon picking challenge," in *Proceedings of the IEEE International Conference on Robotics and Automation*, 2017.
- [4] C. Hernandez, M. Bharatheesha, W. Ko, H. Gaiser, J. Tan, K. van Deurzen, M. de Vries, B. V. Mil, J. van Egmond, R. Burger, M. Morariu, J. Ju, X. Germann, R. Ensing, J. van Frankenhuyzen, and M. Wisse, "Team delft's robot winner of the amazon picking challenge 2016," in *RoboCup 2016: Robot World Cup XX [Leipzig, Germany, June 30 - July 4, 2016]*, ser. Lecture Notes in Computer Science, S. Behnke, R. Sheh, S. Sarel, and D. D. Lee, Eds., vol. 9776. Springer, 2016, pp. 613–624. [Online]. Available: https://doi.org/10.1007/978-3-319-68792-6_51
- [5] E. Shelhamer, J. Long, and T. Darrell, "Fully convolutional networks for semantic segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 4, pp. 640–651, 2017. [Online]. Available: <https://doi.org/10.1109/TPAMI.2016.2572683>
- [6] A. Boularias, J. A. D. Bagnell, and A. T. Stentz, "Learning to manipulate unknown objects in clutter by reinforcement," in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (AAAI-15)*. AAAI, January 2015.
- [7] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *International Conference on Computer Vision (ICCV)*, 2015.
- [8] X. Wang, K. He, and A. Gupta, "Transitive invariance for self-supervised visual representation learning," in *ICCV*, 2017.
- [9] X. Wang and A. Gupta, "Unsupervised learning of visual representations using videos," *CoRR*, vol. abs/1505.00687, 2015. [Online]. Available: <http://arxiv.org/abs/1505.00687>
- [10] X. Wang, K. He, and A. Gupta, "Transitive invariance for self-supervised visual representation learning," 10 2017, pp. 1338–1347.
- [11] S. Purushwalkam and A. Gupta, "Demystifying contrastive self-supervised learning: Invariances, augmentations and dataset biases," 2020.
- [12] J. Mercier, C. Mitash, P. Giguère, and A. Boularias, "Learning object localization and 6d pose estimation from simulation and weakly labeled real images," in *International Conference on Robotics and Automation, ICRA 2019, Montreal, QC, Canada, May 20-24, 2019*. IEEE, 2019, pp. 3500–3506. [Online]. Available: <https://doi.org/10.1109/ICRA.2019.8794112>
- [13] C. Mitash, B. Wen, K. E. Bekris, and A. Boularias, "Scene-level pose estimation for multiple instances of densely packed objects," in *3rd Annual Conference on Robot Learning, CoRL 2019, Osaka, Japan, October 30 - November 1, 2019, Proceedings*, ser. Proceedings of Machine Learning Research, L. P. Kaelbling, D. Kragic, and K. Sugiura, Eds., vol. 100. PMLR, 2019, pp. 1133–1145. [Online]. Available: <http://proceedings.mlr.press/v100/mitash20a.html>
- [14] C. Mitash, A. Boularias, and K. E. Bekris, "Robust 6d object pose estimation with stochastic congruent sets," in *British Machine Vision Conference 2018, BMVC 2018, Newcastle, UK, September 3-6, 2018*. BMVA Press, 2018, p. 277. [Online]. Available: <http://bmvc2018.org/contents/papers/1046.pdf>
- [15] C. Mitash, K. E. Bekris, and A. Boularias, "A self-supervised learning system for object detection using physics simulation and multi-view pose estimation," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2017, Vancouver, BC, Canada, September 24-28, 2017*. IEEE, 2017, pp. 545–551. [Online]. Available: <https://doi.org/10.1109/IROS.2017.8202206>
- [16] C. Mitash, A. Boularias, and K. E. Bekris, "Physics-based scene-level reasoning for object pose estimation in clutter," *CoRR*, vol. abs/1806.10457, 2018. [Online]. Available: <http://arxiv.org/abs/1806.10457>
- [17] —, "Improving 6d pose estimation of objects in clutter via physics-aware monte carlo tree search," in *2018 IEEE International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018*. IEEE, 2018, pp. 1–8. [Online]. Available: <https://doi.org/10.1109/ICRA.2018.8461163>
- [18] I. Misra and L. van der Maaten, "Self-supervised learning of pretext-invariant representations," 2019.
- [19] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," 2020.
- [20] M. Tschannen, J. Djolonga, M. Ritter, A. Mahendran, X. Zhai, N. Houlsby, S. Gelly, and M. Lucic, "Self-supervised learning of video-induced visual invariances," 2020.
- [21] A. Kukleva, H. Kuehne, F. Sener, and J. Gall, "Unsupervised learning of action classes with continuous temporal embedding," 06 2019, pp. 12 058–12 066.
- [22] C. Redondo Cabrera and R. López-Sastre, "Unsupervised learning from videos using temporal coherency deep networks," *Computer Vision and Image Understanding*, 01 2018.
- [23] Y. Yan, H. Hao, B. Xu, J. Zhao, and F. Shen, "Image clustering via deep embedded dimensionality reduction and probability-based triplet loss," *IEEE Transactions on Image Processing*, vol. 29, pp. 5652–5661, 2020.
- [24] G. W. Taylor, R. Fergus, Y. LeCun, and C. Bregler, "Convolutional learning of spatio-temporal features," in *Computer Vision – ECCV 2010*, K. Daniilidis, P. Maragos, and N. Paragios, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 140–153.
- [25] H. Mobahi, R. Collobert, and J. Weston, "Deep learning from temporal coherence in video," in *Proceedings of the 26th Annual International Conference on Machine Learning*, ser. ICML '09. New York, NY, USA: Association for Computing Machinery, 2009, p. 737–744. [Online]. Available: <https://doi.org/10.1145/1553374.1553469>
- [26] D. Stavens and S. Thrun, "Unsupervised learning of invariant features using video," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2010, pp. 1649–1656.
- [27] D. Arthur and S. Vassilvitskii, "K-means++: The advantages of careful seeding," in *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, ser. SODA '07. USA: Society for Industrial and Applied Mathematics, 2007, p. 1027–1035.
- [28] A. Y. Ng, M. I. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," in *Proceedings of the Fourteenth Annual Conference on Neural Information Processing Systems: Natural and Synthetic*, ser. NIPS '01. Cambridge, MA, USA: MIT Press, 2001, p. 849–856.
- [29] S. C. Johnson, "Hierarchical clustering schemes," *Psychometrika*, vol. 32, pp. 241–254, 1967.
- [30] C. Ding and X. He, "K-means clustering via principal component analysis," in *Proceedings of the Twenty-First International Conference on Machine Learning*, ser. ICML '04. New York, NY, USA: Association for Computing Machinery, 2004, p. 29. [Online]. Available: <https://doi.org/10.1145/1015330.1015408>
- [31] R. Luss and A. d'Aspremont, "Clustering and feature selection using sparse principal component analysis," *Optimization and Engineering*, vol. 11, pp. 145–157, 2007.
- [32] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *J. Mach. Learn. Res.*, vol. 11, p. 3371–3408, Dec. 2010.
- [33] J. Xie, R. B. Girshick, and A. Farhadi, "Unsupervised deep embedding for clustering analysis," *CoRR*, vol. abs/1511.06335, 2015. [Online]. Available: <http://arxiv.org/abs/1511.06335>
- [34] B. Yang, X. Fu, N. D. Sidiropoulos, and M. Hong, "Towards k-means-friendly spaces: Simultaneous deep learning and clustering," *CoRR*, vol. abs/1610.04794, 2016. [Online]. Available: <http://arxiv.org/abs/1610.04794>
- [35] C. Song, F. Liu, Y. Huang, L. Wang, and T. Tan, "Auto-encoder based data clustering," in *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, J. Ruiz-Shulcloper and G. Saniti di Baja, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 117–124.
- [36] F. Li, H. Qiao, B. Zhang, and X. Xi, "Discriminatively boosted image clustering with fully convolutional auto-encoders," *CoRR*, vol. abs/1703.07980, 2017. [Online]. Available: <http://arxiv.org/abs/1703.07980>
- [37] K. G. Dizaji, A. Herandi, and H. Huang, "Deep clustering via joint convolutional autoencoder embedding and relative entropy

- minimization,” *CoRR*, vol. abs/1704.06327, 2017. [Online]. Available: <http://arxiv.org/abs/1704.06327>
- [38] P. Huang, Y. Huang, W. Wang, and L. Wang, “Deep embedding network for clustering,” in *2014 22nd International Conference on Pattern Recognition*, 2014, pp. 1532–1537.
 - [39] S. Saito and R. T. Tan, “Neural clustering: Concatenating layers for better projections,” 2017.
 - [40] A. Coates and A. Y. Ng, *Learning Feature Representations with K-Means*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 561–580.
 - [41] C. Hsu and C. Lin, “Cnn-based joint clustering and representation learning with feature drift compensation for large-scale image data,” *CoRR*, vol. abs/1705.07091, 2017. [Online]. Available: <http://arxiv.org/abs/1705.07091>
 - [42] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, “Deep clustering for unsupervised learning of visual features,” *CoRR*, vol. abs/1807.05520, 2018. [Online]. Available: <http://arxiv.org/abs/1807.05520>
 - [43] J. Guérin, O. Gibaru, S. Thiery, and E. Nyiri, “CNN features are also great at unsupervised classification,” *CoRR*, vol. abs/1707.01700, 2017. [Online]. Available: <http://arxiv.org/abs/1707.01700>
 - [44] G. Chen, “Deep learning with nonparametric clustering,” *CoRR*, vol. abs/1501.03084, 2015. [Online]. Available: <http://arxiv.org/abs/1501.03084>
 - [45] V. Lomonaco and D. Maltoni, “Core50: a new dataset and benchmark for continuous object recognition,” ser. Proceedings of Machine Learning Research, S. Levine, V. Vanhoucke, and K. Goldberg, Eds., vol. 78. PMLR, 13–15 Nov 2017, pp. 17–26. [Online]. Available: <http://proceedings.mlr.press/v78/lomonaco17a.html>
 - [46] V. Lomonaco, D. Maltoni, and L. Pellegrini, “Rehearsal-free continual learning over small non-i.i.d. batches,” 2020.
 - [47] K. Lai, L. Bo, X. Ren, and D. Fox, “A large-scale hierarchical multi-view rgb-d object dataset,” in *2011 IEEE International Conference on Robotics and Automation*, 2011, pp. 1817–1824.
 - [48] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, “Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes,” *arXiv preprint arXiv:1711.00199*, 2017.
 - [49] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” *CoRR*, vol. abs/1512.03385, 2015. [Online]. Available: <http://arxiv.org/abs/1512.03385>