

Learning rich touch representations through cross-modal self-supervision

Martina Zambelli
Deepmind, London, UK
zambellim@google.com

Yusuf Aytar
Deepmind, London, UK
yusufaytar@google.com

Francesco Visin
Deepmind, London, UK
visin@google.com

Yuxiang Zhou
Deepmind, London, UK
yuxiangzhou@google.com

Raia Hadsell
Deepmind, London, UK
raia@google.com

Abstract: The sense of touch is fundamental in several manipulation tasks, but rarely used in robot manipulation. In this work we tackle the problem of learning rich touch features from cross-modal self-supervision. We evaluate them identifying objects and their properties in a few-shot classification setting. Two new datasets are introduced using a simulated anthropomorphic robotic hand equipped with tactile sensors on both synthetic and daily life objects. Several self-supervised learning methods are benchmarked on these datasets, by evaluating few-shot classification on unseen objects and poses. Our experiments indicate that cross-modal self-supervision effectively improves touch representation, and in turn has great potential to enhance robot manipulation skills.

Keywords: Touch, Self-supervision, Manipulation

1 Introduction

Deep learning brought us many success stories in a wide range of fields within artificial intelligence, including vision, sound, and text understanding. One paradigm shift that became pervasive recently is that of learning rich high-level feature representations, and then applying them to a wide variety of downstream tasks. These rich representations are learned through supervised pretraining on large scale annotated datasets, such as ImageNet [1, 2], and then applied successfully to learn many downstream tasks with smaller amounts of annotated data, such as object detection [3], segmentation [4] and classification of new categories [5]. This process, termed as *finetuning*, became the main driving force of many new applications of deep learning where annotated data is only scarcely available.

Most recently, the self-supervised learning approaches, which aim to learn rich representations from uncurated large scale data, removed the costly burden of annotations and enabled learning rich high-level representations in a large range of domains. In the image understanding domain self-supervised representations already reached the quality of representations learned through large scale annotated datasets [6, 7]. Self-supervised pretraining became the de-facto standard in text understanding [8, 9, 10] and even modalities like sound greatly benefited from rich representations learned through self-supervised approaches utilizing multiple modalities [11, 12, 13]. This methodology of obtaining rich features that enable *sample-efficient* (i.e. few-shot) learning strongly resonates with us, as humans mostly learn new tasks in a similar manner even with very small amounts of data or experiences. In this paper we propose to learn rich touch representations that enable sample-efficient learning for touch-based perception (i.e. perception in the dark). We achieve this through cross-modal self-supervised learning by harnessing the natural correlation between vision and touch signals.

The sense of touch is ubiquitous in our everyday life. While vision is possibly the predominant sensory modality through which we learn about and interact with our surroundings, touch plays a critical role in manipulation tasks. Indeed, it provides us with very heterogeneous and rich data, with direct measurements of ongoing contact forces during object interactions that allow us to infer friction, compliance, mass, and other physical properties of surfaces and objects. Critically, the mental models built combining data from vision, touch and possibly other modalities are very effectively exploited when vision is not available.

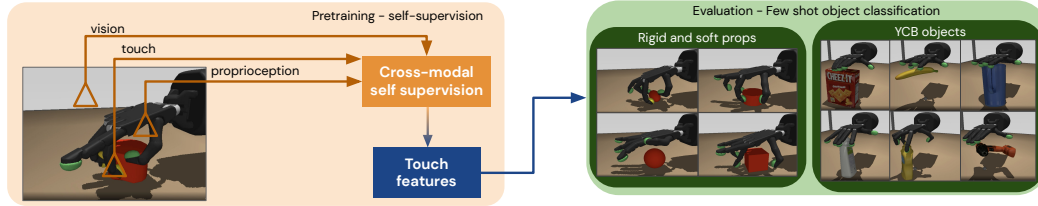


Figure 1: Learning touch features through cross-modal self-supervision and evaluation on few-shot object classification. Learned features can in turn be used to identify different objects and their properties across diverse sets of synthetic and real-world objects.

As a motivating example, consider a very common task such as looking for a set of keys in a bag. While we typically have a clear picture of how the keys look like, it is often hard to *look* properly inside a bag, and we have to resort to other senses. In this situation, reaching with a hand in the bag allows us to explore using touch alone but, critically, exploiting the mental model previously built using data from the richer vision modality. This example is representative of a broader set of tasks where visual clues are limited or precluded when the task is performed, but are usually accessible otherwise. Consider for instance shifting gears while keeping the eyes on the road, grasping a glass of water in the dark or tying shoelaces without watching. While these tasks are relying mainly, if not only, on touch, they exploit a mental model of the objects that has been built throughout the person’s life using all senses.

Such tasks are very common in our everyday life, and critically also in robot manipulation. Here, occlusions, restrictions of the visual field, or limited vision capabilities, can be a real impediment for the successful completion of many tasks. For instance understanding the identity and orientation of grasped objects when occluded by a robot hand can be challenging relying only on vision, but becomes much easier when paired with touch and proprioception. Indeed, vision alone is often not sufficient to properly execute dexterous manipulation tasks. A representation learned using multiple sensor modalities, e.g. vision, touch and proprioception, can be fundamental to capture key features of objects.

In this paper, we present an extensive study of cross-modal techniques that learn rich representations from vision and touch, without requiring external supervision. We apply many existing multi-modal representation learning methods, such as L^3 -net [11], CMDC [14], TCN [15] and CMC [16], to vision and touch modalities. We develop two cross-modal versions of CPC [17] and implement a cross-modal generation method (CM-GEN), which generates the visual observations using the touch features. We evaluate the effectiveness of learned touch representations on few-shot classification of objects and their orientations with increasing levels of complexity. A visual representation of our setup is shown in Fig. 1. We demonstrate that learned touch features are significantly better than using raw proprioception and touch features, particularly in few-shot learning scenarios where we obtain a 25% average performance boost. We release two datasets¹ with the rich multimodal data we trained on, to enable further research in cross-modal representation learning with touch and vision. These are based on a simulated Shadow Dexterous hand [18], and include visual, tactile and proprioceptive trajectories.

2 Related work

When it comes to robotics, the importance of tactile sensing has long been acknowledged [19, 20, 21, 22] but making use of it has proven challenging for a number of reasons, from hardware related issues like wear and tear and drift, to the lack of simple ways to specify tactile-based goals or trajectories, as well as to the limitations of physics models of contacts and forces. Recent works have explored methods to combine vision and touch to learn representations from robot interactions with objects [23, 24, 25, 26, 27] and several methods have been proposed to learn multi- and cross-modal representations from large datasets [28, 29, 30, 31, 32, 33, 34, 35] in other machine learning domains, but a comprehensive evaluation of how complete and useful these representations are on practical tasks is still lacking.

Cross-modal learning Our work is mostly related to cross-modal representation learning [28, 29], which aims to generate a representation that effectively correlates different sensor modalities. Several studies explored cross-modal transfer between images and text [36], as well as models for vision and language representations [30, 31]. Interesting work has been done on audio-visual representations [32, 33, 34, 35] and on generating captions for images [37, 38, 39, 40]. Large-scale cross-domain datasets

¹accessible at <https://github.com/deepmind/deepmind-research/tree/master/cmtouch>

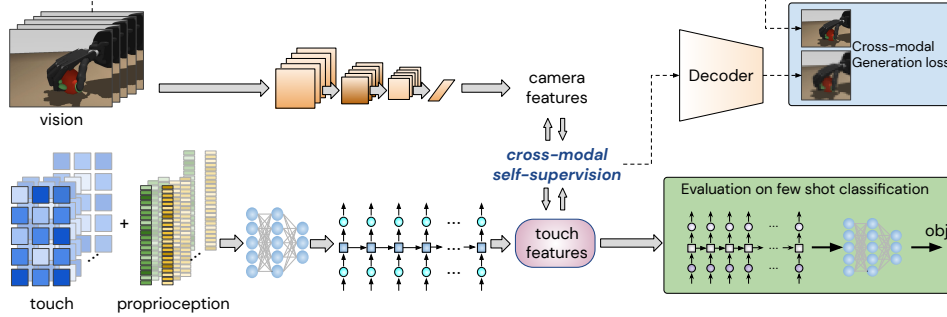


Figure 2: Overview: self-supervised cross-modal representation learning of touch features. Cross-modal generation loss used to train features (top-right), and evaluation networks (bottom-right).

were fundamental to their success, but these types of datasets are currently unavailable for vision and touch. We hence used a robotic platform equipped with tactile sensors to collect data for our study.

We are specifically interested in the interplay between vision and touch. Recent studies have addressed the problem of integrating these sensor modalities by studying the physical properties of fabrics via visual and tactile data fusion [24, 25, 23, 41], or by focusing on multimodal learning using proprioception, vision, touch and sound [42, 27, 43, 44]. Different from prior work that addressed specific manipulation tasks, we focus on several cross-modal representation methods to learn touch features through self-supervision. Rather than shaping the learned representation around a specific downstream manipulation task (such as picking-up [45], grasping [44], classification [46, 23]), the proposed approach focuses on learning a representation independent of the target use case using self-supervision.

3 Learning Rich Touch Features

In this section we describe how high-level touch features are learned by exploiting naturally existing correlations between touch perception and vision. No annotations or supervision is provided for feature learning, i.e., the process is entirely *self-supervised*. As described in Section 4, our episodic datasets of robotic hand and object interactions are obtained by running a policy trained blindly to explore its surroundings. The features are learned by capitalizing on temporal correlations of vision and touch within the provided episodes (Fig. 2).

Formally, for any given paired vision $\mathbf{o} = \{o_i\}_{i=1}^N$ and touch $\mathbf{t} = \{t_i\}_{i=1}^N$ sequences, where o_i and t_i are visual and touch observations at time step i , we first encode them to obtain the vision and touch features $v_i = f(o_i)$ and $x_i = g(t_i)$ respectively, and then we employ several cross-modal self-supervision objectives to learn latent touch and vision features. Next we’ll present the touch and vision encoders, then we’ll describe the cross-modal self-supervision objectives used for learning touch features.

3.1 Encoders

Touch Encoding. Touch and proprioception are encoded through a 4-layers MLP followed by a LSTM network, which allows to accumulate the dynamics information provided by tactile sensors and robot joints throughout an episode. Capturing this information in a time window is critical to learn meaningful features that reflect the dynamical nature of these signals.

Vision Encoding. Visual input is obtained from simulated cameras, and is processed through a convolutional encoder with residual connections. Each frame is processed separately. Contrary to proprioception and touch information, visual information is not accumulated with any memory unit.

Implementation details. All implementation details are reported in the Supplementary material.

3.2 Cross-modal self-supervision

We adapted several cross-modal self-supervision objectives, such as L^3 -net[11], CMDC[14], TCN[15], and CMC[16] to learn touch and vision features. We developed cross-modal versions of contrastive predictive coding (CPC)[17] with two variants. We also implemented a cross-modal generation method CM-GEN which generates the visual observations using the touch features. All the methods operate over the latent paired sequences $\mathbf{v} = \{v_i\}_{i=1}^N$ and $\mathbf{x} = \{x_i\}_{i=1}^N$ except for the CM-

GEN which also uses o_i while generating visual frames. We describe each of these methods below using a single paired vision and touch sequence, but we optimize them with all the episodes in our dataset.

L³-net[11]. The main idea in L³-net is binary classification of temporally sync and non-sync pairs in multiple modalities. In our case a positive example is a synced pair (v_i, x_i) and negative examples are any non-sync pairs $(v_i, x_j), \forall i \neq j$. Non-sync pairs are also acquired by mixing vision and touch latents from multiple episodes. Given a positive or a negative pair, we concatenate touch and vision features (v, x) and run them through a two-level MLP to perform binary sync/non-sync classification. The loss function is binary cross-entropy.

CMDC[14]. Instead of directly using synced pairs, cross-modal distance classification (CMDC) predicts the temporal distance between the given cross-modal pair. For instance in the pair (v_i, x_{i+1}) the temporal distance is 1. CMDC sets distance prediction as a classification problem where the classes are different distance intervals. We follow the same setting described in [14]. Given (v, x) we run them through a two-layers MLP to perform distance classification with a softmax cross-entropy loss.

TCN[15] and CMC[16]. Time contrastive networks (TCN) learn representations by enforcing higher similarity between sync pairs across two different camera views compared to any pair selected within a single view sequence. In our context two views are the vision and touch modalities. We particularly adopted n-pairs implementation [15] as opposed to the triplet implementation as the model learns to cheat in triplet implementation when there is an LSTM in one of the modalities (i.e. touch). Given the paired sequences \mathbf{v} and \mathbf{x} the TCN n -pairs objective that we minimize is:

$$-\log \frac{\exp(x_i^T v_i)}{\sum_j \exp(x_i^T v_j)} \quad (1)$$

where (x_i, v_i) is the positive pair and $(x_i, v_j), \forall j \neq i$ are the negative pairs. This specific implementation has very strong connection to contrastive multi-view coding (CMC[16]) and multi-modal noise-contrastive estimation (NCE[12]) which are essentially the same objectives. The aim in both methods is to classify the matching pair among all other non-matching negative pairs. The main difference compared to TCN is that negative pairs are also populated from other sequences in the batch. In equation (1) we can also swap vision and touch to obtain the symmetric loss function. In our implementations of TCN and CMC we use the symmetric losses as well.

CPC[17]. Contrastive predictive coding (CPC) predicts the future latent observations using the current latent observation and previous context which is accumulated through an LSTM. It is originally developed for single modal sequences. Here we build a cross-modal version where a touch feature x_i predicts n -step into the future in the vision modality. Each of the n -step predictions are obtained through n linear predictors for each time step distance. The objective is formalized as a contrastive learning approach where (v_i, \hat{v}_i) is the positive pair and \hat{v}_i is the predicted vision feature using an earlier touch feature. Note that we have n different predictions of v_i as any previous $x_{i-1}, x_{i-2}, \dots, x_{i-n}$ can provide their own prediction \hat{v}_i . The minimized objective is:

$$-\log \frac{\exp(\hat{v}_i^T v_i)}{\sum_j \exp(\hat{v}_i^T v_j)} \quad (2)$$

where we contrast the matching prediction with all other non-matching predictions. We refer to this method as cross-modal CPC (CM-CPC) as it is a straightforward adaptation of CPC to multiple modalities. Unlike the regular CPC, as we operate in multiple modalities we can also predict the current or previous latents in the other modality. Hence we also introduce CM-CPC-H, a variation of CM-CPC that predicts the n -step history, current, and n -step future latents in the vision modality. In practice we set $n = 20$, similar to [17]. For more details on CPC please refer to [17].

CM-GEN. Generating one modality from another is performed numerous times in multiple domains for different purposes [47, 48, 49, 50]. In our context we generate vision from touch mainly for learning high-level touch representations. We employ a 5 layers convolutional decoder d to the touch feature x_i to generate the corresponding visual frame o_i . The minimized loss is $\|d(x_i) - o_i\|^2$ for all the time steps i . The details of the decoder can be found in the Supplementary material.

Implementation details. We use Adam [51], with learning rate 0.001 for CM-GEN and 0.0001 for all other methods. All the other implementation details are reported in the Supplementary material.

4 Dataset and Experimental Setup

The proposed approach allows us to learn touch features through cross-modal self-supervision (Fig. 2). The learned features are then evaluated on few-shot classification tasks, to identify different objects and their properties. The following paragraphs describe how the datasets have been generated.

Experimental setup We run experiments in simulation with MuJoCo [52] and we use the simulated Shadow Dexterous Hand [18], with five fingers and 24 degrees of freedom, actuated by 20 motors. In simulation, each fingertip has a spatial touch sensor attached with a spatial resolution of 4×4 and three channels: one for normal force and two for tangential forces. We simplify this by taking the absolute value and then summing across the spatial dimensions, to obtain a 3D force vector for each fingertip. More details about the tactile sensors are reported in the Supplementary material.

The state consists of proprioception (joint positions and joint velocities) and touch. Visual inputs are collected with a 64×64 resolution and are only used for representation learning, but are not provided as observations to control the robot’s actions. The action space is 20-dimensional. We use velocity control and a control rate of 30 Hz. Each episode has 200 time steps, which correspond to about 6 seconds.

The environment consists of the Shadow Hand, facing down, and interacting with different objects. These objects have different shapes, sizes and physical properties (e.g. rigid or soft). We develop two versions of the task, the first using simple props and the second using YCB objects. In both cases, objects are fixed to their frame of reference, while their position and orientation are randomized. Focusing on this simple environment enables us to clearly characterize the effectiveness of our methods for training cross-modal and touch-based representations.

Props Props are simple 3D shaped objects that include cubes, spheres, cylinders and ellipsoid of different sizes. We also generated the soft version of each prop, which can deform under the pressure of the touching fingers. Soft deformable objects are complex entities to simulate: they are defined through a composition of multiple bodies (capsules) that are tied together to form a shape, such as a cube or a sphere. The main characteristic of these objects is their elastic behaviour, that is they change shape when touched. The most difficult thing to simulate in this context is contacts, which grow exponentially with the increased number of colliding bodies. Forty-eight different objects are generated by sampling from 6 different sized, 4 different shapes (i.e. sphere, cylinder, cube, ellipsoid), and they can either be rigid or soft. A sample from this dataset is shown in Fig. 3.

YCB objects The YCB objects dataset [53] consists of everyday objects with different shapes, sizes, textures, weight and rigidity. We chose a set of ten objects: cracker box, sugar box, mustard bottle, potted meat can, banana, pitcher base, bleach cleanser, mug, power drill, scissors. These are generated in simulation at their standard size, which is also proportionate to the default dimension of the simulated Shadow Hand (see examples in Fig. 4). The pose of each object is randomly selected among a set of 60 different poses, where we vary the orientation of the object. These variations make the identification of each object more complex and require a higher generalization capability from the learning method applied. More details about the dataset generation are reported in the Supplementary material. A sample from this dataset is shown in Fig. 4.

Control A Maximum a Posteriori Policy Optimization (MPO) [54] reinforcement learning (RL) agent is used to control the hand movements. At the beginning of each episode, an object is chosen randomly among the set of objects. The agent is rewarded when touching the object with the sensorized fingertips: the more fingertips are in contact with the object, the higher the reward. The agent is blind, as no visual input is fed to its networks, and controls the hand by issuing velocity commands. It is important that no visual information is provided to the policy: when vision is present as input to the agent, visual clues are exploited to generate control trajectories. In turn, the data collected, specifically proprioception, would contain significant biases leaked from the visual clues given to the policy. By training a blind agent we make sure that proprioception and touch are not biased by any visual clue.

5 Experimental Results

In this section we include quantitative results of few-shot classification, evaluation on unseen objects and pose estimation. A discussion follows.

Few-shot classification We evaluate all our training methods on few-shot classification from the learned touch features. Vision is not available at this time, just as when looking for the keys in a bag.

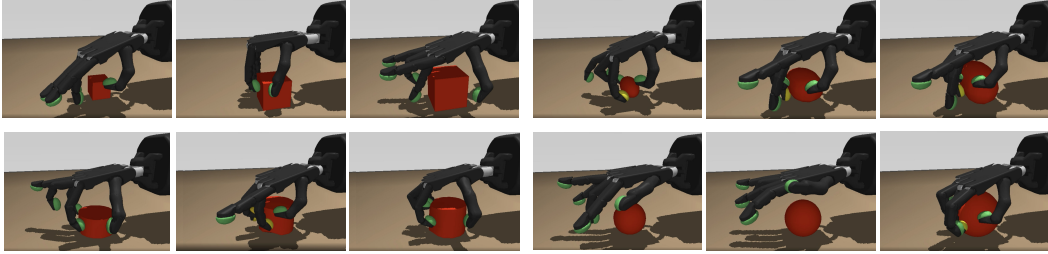


Figure 3: Samples from the Props dataset: the Shadow hand is presented with objects sampled from four shapes of different sizes, and they can be either rigid or deformable objects.

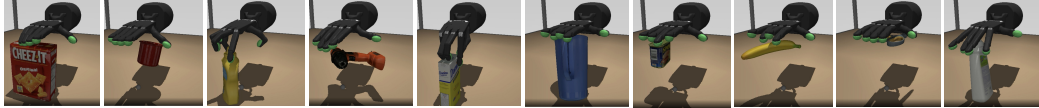


Figure 4: Samples from the YCB objects dataset: the Shadow hand is presented with objects sampled from ten classes, which appear in different orientations and poses.

Table 1: Few-shot classification - Props and YCB dataset.

	Props: all → all								YCB: all → all							
	1	3	5	10	25	50	100	mean	1	3	5	10	25	50	100	mean
CM-CPC	11.3%	23.4%	29.8%	52.0%	70.7%	84.0%	88.2%	51.3%	14.4%	20.7%	19.3%	29.0%	35.2%	44.9%	55.9%	31.4%
CM-CPC-H	12.0%	24.2%	31.1%	52.3%	74.2%	83.2%	86.7%	52.0%	16.2%	20.7%	20.8%	27.8%	34.3%	49.3%	55.4%	32.1%
CM-GEN	18.1%	42.6%	54.0%	72.6%	79.3%	86.7%	88.6%	63.1%	27.6%	41.0%	49.7%	50.9%	58.1%	66.6%	61.2%	50.7%
CMC	11.9%	24.6%	33.0%	49.6%	68.9%	79.3%	84.0%	50.2%	14.5%	17.0%	21.1%	27.3%	35.2%	44.7%	57.2%	31.0%
CMDC	11.9%	24.9%	33.3%	53.9%	71.5%	82.6%	87.2%	52.2%	15.0%	17.5%	18.9%	28.0%	36.2%	48.5%	58.8%	31.8%
L3-NET	11.4%	24.1%	28.4%	49.3%	69.3%	80.6%	87.4%	50.1%	16.3%	16.6%	19.6%	24.6%	32.0%	46.9%	58.5%	30.6%
TCN	13.2%	25.8%	32.4%	53.1%	68.9%	77.8%	83.0%	50.6%	19.8%	15.8%	21.3%	23.3%	36.1%	45.9%	57.3%	31.4%
scratch	09.2%	23.9%	28.1%	48.6%	64.5%	79.2%	84.1%	48.2%	13.1%	15.2%	17.9%	25.0%	30.1%	36.0%	48.7%	26.6%

We compare our self-supervised approach with a learning from scratch approach, where features are learned directly from touch and proprioception while optimizing for the classification objective. For the Props dataset, a class is defined by an instance of a specific shape (cube, sphere, cylinder or ellipsoid), size (one of 6 different sizes) and rigidity (rigid or soft object). There is therefore a total of 48 classes for the props dataset. There are instead 10 classes for the YCB objects dataset, corresponding to the 10 different objects used to collect the data. Classification of the 10 YCB objects is nonetheless challenging due to the variations applied to the pose of each object in each episode.

Results of few-shot classification are presented in Table 1. Learned features are learned on the entire datasets, and are evaluated on few-shot classification using 1, 3, 5, 10, 25, 50 or 100 examples for each object class. The mean across scores is also reported. Darker colors indicate better scores. The standard deviation obtained across multiple runs of the experiment is between 0.01 and 0.04.

Evaluation on unseen objects We evaluate learned touch features on objects that were not presented in the pretraining step. In order to do so, we generate two disjoint sets of objects from the collected datasets, that we call simply Set1 and Set2. One of the sets is used to train the self-supervised representation, and the other set is used to evaluate the learned touch features. The two sets are randomly generated to be disjoint, that is Set1 and Set2 don't have any object in common, hence models are evaluated on completely unseen objects. Results of evaluation on unseen objects from Set1 to Set2 and from Set2 to Set1 are presented in Figures 5 and 6 for the Props and YCB objects datasets, respectively. Similar to the previous experiment, the standard deviation obtained across multiple runs of this experiment is between 0.01 and 0.04. Further results are reported in the Supplementary material, including all accuracy scores of the different methods.

We also perform a cross-dataset evaluation: touch features learned on the props datasets are evaluated on few-shot classification of YCB objects, and viceversa. Results of this evaluation are presented in Fig. 7. For both sets of experiments, learned features are trained on one of the datasets, and evaluated

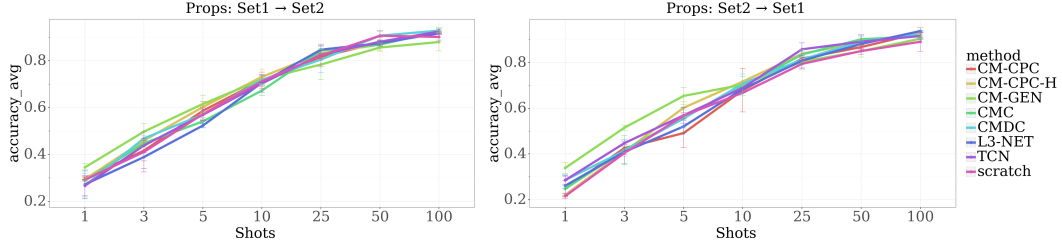


Figure 5: Evaluation on unseen objects - Props objects classification. In the low shots regime CM-GEN outperforms other methods, but above 10 shots the difference is less prominent, with CMDC and TCN showing slightly better performance than other methods in average.

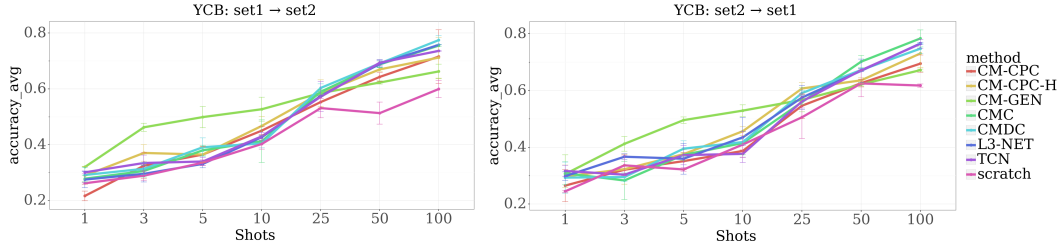


Figure 6: Evaluation on unseen objects - YCB objects classification. Similar to the results on the Props dataset, in the low shots regime CM-GEN outperforms other methods. Above the 25 shots regime, CMC, CMDC, TCN and L3-NET tend to perform best with similar performance.

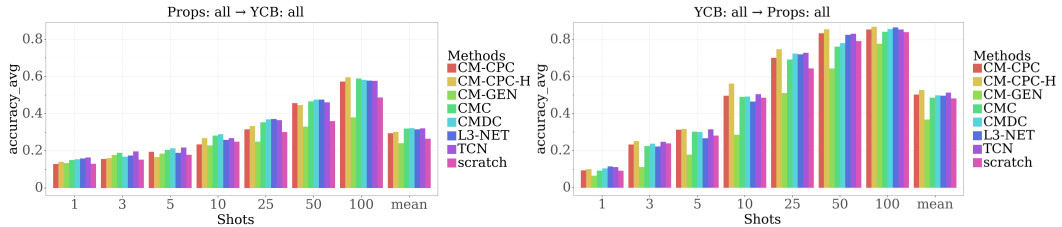


Figure 7: Evaluation on unseen objects across datasets.

on few-shot classification on the other dataset, using 1, 3, 5, 10, 25, 50 or 100 examples for each object’s class.

Pose estimation For the YCB objects dataset, we also evaluated the learned touch features on estimating the orientation of each object class. Pose estimation is formulated as a classification across a set of different poses that an object can take: each object can appear in 60 different poses, given by a rotation around the vertical axis and a tilt from this same axis.

Results of few-shot pose estimation are presented in Fig. 8, where for each object we report the accuracy obtained with only 1, 10, 50, 100 examples of that object. Table 2 reports the average accuracy across all objects, for all methods and all few-shots experiments. Further results are reported in the Supplementary material, including all few-shot (1, 3, 5, 10, 25, 50 or 100) experiments.

5.1 Discussion

Results presented above show that cross-modal self-supervision improves touch features: touch features learned from scratch achieved lower scores compared to features learned through our proposed cross-modal self-supervision approach.

In all experiments, as expected, we notice an increase in performance when more examples are provided for n -shot classification (i.e., when n is higher). The accuracy is low when only a single example is provided (18.1% for props and 27.6% for YCB objects - Table 1), and grows with the number of seen examples, reaching an accuracy of 88.6% and 61.2% for props and YCB objects, respectively, when 100 examples are available. Note that props are generally easier to identify compared to YCB objects, given the simple shapes and the different textures (rigid or soft).

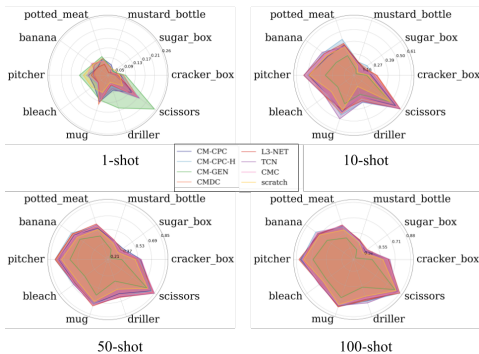


Figure 8: Pose estimation accuracy of YCB objects: 1-shot, 10-shot, 50-shot and 100-shot classification. Larger area is better.

	1	3	5	10	25	50	100
CM-CPC	08.7%	17.3%	23.4%	35.1%	51.0%	61.4%	68.1%
CM-CPC-H	08.7%	16.7%	24.4%	38.3%	52.9%	62.6%	68.1%
CM-GEN	11.6%	19.8%	23.1%	26.5%	36.7%	48.3%	56.1%
CMC	08.5%	16.3%	23.4%	35.8%	52.7%	62.0%	67.6%
CMDC	08.5%	15.6%	23.4%	36.7%	51.0%	62.5%	68.2%
L3-NET	07.6%	16.0%	23.4%	37.6%	52.4%	63.2%	68.5%
TCN	08.2%	16.7%	23.9%	37.8%	53.5%	62.8%	68.5%
scratch	07.9%	14.6%	20.9%	31.4%	45.0%	56.9%	64.4%

Table 2: Pose estimation: accuracy across all objects. Learned features are evaluated on few-shot classification using 1, 3, 5, 10, 25, 50 or 100 examples for each object’s class. Darker colors indicate better scores.

For few-shot classification evaluated on single dataset, CM-GEN performed best compared to other self-supervision methods. This is because visual clues captured during self-supervision are sufficient and effective. Other methods become competitive when training is performed on a set of objects and the learned features are then evaluated on a second set. In this case, results confirm that cross-modal self-supervision achieves considerably better performance than training touch features from scratch. CM-GEN still achieves slightly better results in the low data regime (when only a small number of classes are observed), while other methods based on contrastive approaches (e.g. TCN and L^3 -NET) achieve better results with 50 or 100 examples (Figures 5 and 6). We also note that CM-GEN overfits to the dataset it is trained on, achieving the worst performance in the identification of objects that belong to a different dataset (Fig. 7). In this case, transferring from one group of objects to the other requires a higher level of abstraction, which does not necessarily favour visual clues. These results indicate that methods based on contrastive self-supervised losses, such as TCN, L^3 -NET, CMDC and CMC can learn more robust touch features compared to CM-GEN and thus achieve the best overall performance.

Finally, the proposed approach achieves good performance also on the pose estimation experiment, where the classification is harder due to the large number of classes per object. As observed above, cross-modal generation performs poorly on this task, due to the increased complexity: discriminating orientation from vision is hard and doesn’t provide features that are robust enough. On the other hand, all other self-supervised methods achieve better performance compared to learning features from scratch, thus showing that the proposed approach is an effective way to improve touch representation for manipulation skills with different objects.

Validation on real data While no real world data have been used in these experiments, the proposed approach is agnostic to the specific kind of data used: no bias is introduced in the data format, and the encoder networks used to train the self-supervised representation can be easily adapted, if needed, to accommodate inputs from a real robotic platform. More details in the Supplementary material.

6 Conclusion

The use of touch is bound to have a critical role in robotics manipulation, especially in tasks involving dexterity, manipulation of daily-life objects, dealing with occlusions or closed and constrained environment. We took inspiration from a challenging task: finding an object in a bag. In this paper we present an approach to use touch and vision to build representations that allow to identify objects from touch and proprioception alone. We benchmark several self-supervised learning methods on three challenging tasks: few-shot classification, knowledge transfer to unseen objects, and pose estimation from touch. Experimental results show that cross-modal self-supervision enables the learning of rich touch features that are effective on the final task. In fact, our evaluations demonstrate that touch features learned through cross-modal self-supervised approaches improve classification performance on few-shot and unseen objects classification, as well as pose estimation of several real-world objects.

Furthermore, to spark interest in these types of tasks and to foster further research on the use of touch in robotic manipulation, we release two rich multimodal datasets used in this work.

References

- [1] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [3] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Region-based convolutional networks for accurate object detection and segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 38(1):142–158, 2015.
- [4] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
- [5] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *International conference on machine learning*, pages 647–655, 2014.
- [6] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020.
- [7] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.
- [8] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever. Improving language understanding by generative pre-training, 2018.
- [9] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [10] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5753–5763, 2019.
- [11] R. Arandjelovic and A. Zisserman. Look, listen and learn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 609–617, 2017.
- [12] J.-B. Alayrac, A. Rezasens, R. Schneider, R. Arandjelović, J. Ramapuram, J. De Fauw, L. Smaira, S. Dieleman, and A. Zisserman. Self-supervised multimodal versatile networks. *arXiv preprint arXiv:2006.16228*, 2020.
- [13] B. Korbar, D. Tran, and L. Torresani. Cooperative learning of audio and video models from self-supervised synchronization. In *Advances in Neural Information Processing Systems*, pages 7763–7774, 2018.
- [14] Y. Aytar, T. Pfaff, D. Budden, T. Paine, Z. Wang, and N. de Freitas. Playing hard exploration games by watching youtube. In *Advances in Neural Information Processing Systems*, pages 2930–2941, 2018.
- [15] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, S. Levine, and G. Brain. Time-contrastive networks: Self-supervised learning from video. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1134–1141. IEEE, 2018.
- [16] Y. Tian, D. Krishnan, and P. Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019.
- [17] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [18] ShadowRobot. ShadowRobot Dexterous Hand. <https://www.shadowrobot.com/products/dexterous-hand/>.
- [19] H. Yousef, M. Boukallel, and K. Althoefer. Tactile sensing for dexterous in-hand manipulation in robotics a review. *Sensors and Actuators A: physical*, 167(2):171–187, 2011.
- [20] R. S. Dahiya, G. Metta, M. Valle, and G. Sandini. Tactile sensing from humans to humanoids. *IEEE transactions on robotics*, 26(1):1–20, 2009.
- [21] M. L. Hammock, A. Chortos, B. C.-K. Tee, J. B.-H. Tok, and Z. Bao. 25th anniversary article: the evolution of electronic skin (e-skin): a brief history, design considerations, and recent progress. *Advanced materials*, 25(42):5997–6038, 2013.

- [22] L. Cramphorn, B. Ward-Cherrier, and N. F. Lepora. Tactile manipulation with biomimetic active touch. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 123–129. IEEE, 2016.
- [23] J.-T. Lee, D. Bollegala, and S. Luo. “touching to see” and “seeing to feel”: Robotic cross-modal sensory data generation for visual-tactile perception. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 4276–4282. IEEE, 2019.
- [24] W. Yuan, S. Wang, S. Dong, and E. Adelson. Connecting look and feel: Associating the visual and tactile properties of physical materials. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5580–5588, 2017.
- [25] W. Zheng, H. Liu, and F. Sun. Lifelong visual-tactile cross-modal learning for robotic material perception. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [26] J. Lin, R. Calandra, and S. Levine. Learning to identify object instances by touch: Tactile recognition via multimodal matching. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 3644–3650. IEEE, 2019.
- [27] M. Zambelli, A. Cully, and Y. Demiris. Multimodal representation models for prediction and control from partial information. *Robotics and Autonomous Systems*, 123:103312, 2020.
- [28] Y. Aytar, L. Castrejon, C. Vondrick, H. Pirsivash, and A. Torralba. Cross-modal scene networks. *IEEE transactions on pattern analysis and machine intelligence*, 40(10):2303–2314, 2017.
- [29] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng. Multimodal deep learning. In *ICML*, 2011.
- [30] J. Lin, A. Yang, Y. Zhang, J. Liu, J. Zhou, and H. Yang. Interbert: Vision-and-language interaction for multi-modal pretraining. *arXiv preprint arXiv:2003.13198*, 2020.
- [31] J. Lu, D. Batra, D. Parikh, and S. Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, pages 13–23, 2019.
- [32] N. Le and J.-M. Odobez. Improving speech embedding using crossmodal transfer learning with audio-visual data. *Multimedia Tools and Applications*, 78(11):15681–15704, 2019.
- [33] A. Owens, P. Isola, J. McDermott, A. Torralba, E. H. Adelson, and W. T. Freeman. Visually indicated sounds. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2405–2413, 2016.
- [34] A. Owens, J. Wu, J. H. McDermott, W. T. Freeman, and A. Torralba. Ambient sound provides supervision for visual learning. In *European conference on computer vision*, pages 801–816. Springer, 2016.
- [35] Y. Aytar, C. Vondrick, and A. Torralba. Soundnet: Learning sound representations from unlabeled video. In *Advances in neural information processing systems*, pages 892–900, 2016.
- [36] R. Socher, M. Ganjoo, C. D. Manning, and A. Ng. Zero-shot learning through cross-modal transfer. In *Advances in neural information processing systems*, pages 935–943, 2013.
- [37] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2625–2634, 2015.
- [38] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137, 2015.
- [39] G. Kulkarni, V. Premraj, V. Ordonez, S. Dhar, S. Li, Y. Choi, A. C. Berg, and T. L. Berg. Babytalk: Understanding and generating simple image descriptions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(12):2891–2903, 2013.
- [40] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057, 2015.
- [41] M. A. Lee, Y. Zhu, K. Srinivasan, P. Shah, S. Savarese, L. Fei-Fei, A. Garg, and J. Bohg. Making sense of vision and touch: Self-supervised learning of multimodal representations for contact-rich tasks. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8943–8950. IEEE, 2019.
- [42] H. Liu, F. Sun, B. Fang, and D. Guo. Cross-modal zero-shot-learning for tactile object recognition. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 2018.

- [43] M. Zambelli and Y. Demiris. Online multimodal ensemble learning using self-learned sensorimotor representations. *IEEE Transactions on Cognitive and Developmental Systems*, 9(2):113–126, 2016.
- [44] R. Calandra, A. Owens, M. Upadhyaya, W. Yuan, J. Lin, E. H. Adelson, and S. Levine. The feeling of success: Does touch sensing help predict grasp outcomes? *arXiv preprint arXiv:1710.05512*, 2017.
- [45] Y. Chan, H. Yu, and R. P. Khurshid. Effects of force-torque and tactile haptic modalities on classifying the success of robot manipulation tasks. In *2019 IEEE World Haptics Conference (WHC)*, pages 586–591. IEEE, 2019.
- [46] P. Falco, S. Lu, C. Natale, S. Pirozzi, and D. Lee. A transfer learning approach to cross-modal object recognition: From visual observation to robotic haptic exploration. *IEEE Transactions on Robotics*, 35(4): 987–998, 2019.
- [47] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.
- [48] L. Chen, S. Srivastava, Z. Duan, and C. Xu. Deep cross-modal audio-visual generation. In *Proceedings of the on Thematic Workshops of ACM Multimedia 2017*, pages 349–357, 2017.
- [49] Q. Chen and V. Koltun. Photographic image synthesis with cascaded refinement networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1511–1520, 2017.
- [50] A. Salvador, M. Drozdal, X. Giro-i Nieto, and A. Romero. Inverse cooking: Recipe generation from food images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10453–10462, 2019.
- [51] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2017.
- [52] E. Todorov, T. Erez, and Y. Tassa. MuJoCo: A physics engine for model-based control. In *Proceedings of the International Conference on Intelligent Robots and Systems (IROS)*, 2012.
- [53] B. Calli, A. Singh, J. Bruce, A. Walsman, K. Konolige, S. Srinivasa, P. Abbeel, and A. M. Dollar. Yale-cmu-berkeley dataset for robotic manipulation research. *The International Journal of Robotics Research*, 36(3):261–268, 2017.
- [54] A. Abdolmaleki, J. T. Springenberg, Y. Tassa, R. Munos, N. Heess, and M. Riedmiller. Maximum a posteriori policy optimisation. *arXiv preprint arXiv:1806.06920*, 2018.