

Overview Paper

Bridging Gap between Image Pixels and Semantics via Supervision: A Survey

Jiali Duan and C.-C. Jay Kuo*

University of Southern California, Los Angeles, CA, USA

ABSTRACT

The fact that there exists a gap between low-level features and semantic meanings of images, called the semantic gap, is known for decades. Resolution of the semantic gap is a long standing problem. The semantic gap problem is reviewed and a survey on recent efforts in bridging the gap is made in this work. Most importantly, we claim that the semantic gap is primarily bridged through supervised learning today. Experiences are drawn from two application domains to illustrate this point: (1) object detection and (2) metric learning for content-based image retrieval (CBIR). To begin with, this paper offers a historical retrospective on supervision, makes a gradual transition to the modern data-driven methodology and introduces commonly used datasets. Then, it summarizes various supervision methods to bridge the semantic gap in the context of object detection and metric learning.

Keywords: Semantic Gap, Semantic Understanding, Content-based Image Retrieval, Supervision, Object Detection, Metric Learning.

1 Introduction

Computer vision deals with how computers can gain high-level understanding of visual contents, which are represented by pixels. High-level understanding of visual inputs demands the capability to learn the semantics conveyed through

*Corresponding author: Jiali Duan, jialidua@usc.edu

raw pixels. The fact that there exists a gap between low-level features and semantic meanings of images, is known for decades. It is the consequence of “lack of coincidence between the information that one can extract from the visual data and the interpretation that the same data have for a user in a given situation” [139]. It is well known as the semantic gap.

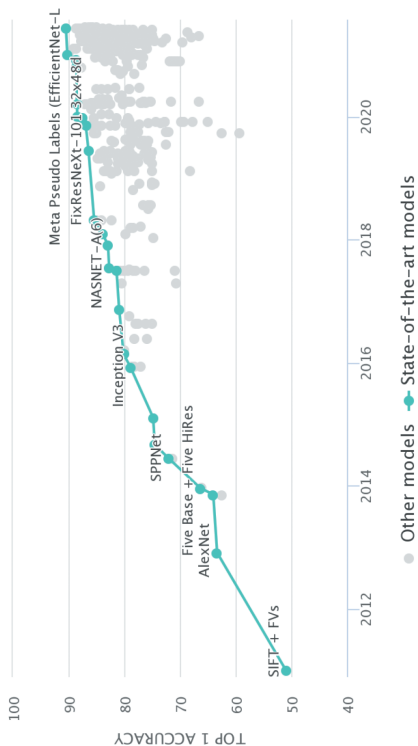
There is a growing awareness in the computer vision community that the key to today’s vision problems lies in resolving the gap between image pixels and semantics. There has been a substantial amount of progress in bridging the gap in recent years. In our opinion, this advancement is primarily attributed to supervised learning and our survey paper is written around this central idea. Supervision manifests itself through two aspects: (1) large-scale, high-quality annotated data, and (2) well-designed optimization objectives. The two aspects often come into play synergistically. For example, the design of optimization objectives highly depend on annotations. The optimization procedure often entails a minimum amount of labeled data and it is expected to scale well with more data.

Table 1 gives an example of how supervision in the form of “annotated data” advances the object classification field. The figure shows the progress of top-1 classification accuracy as a function of time with respect to the ImageNet dataset. Undoubtedly, the large-scale annotated ImageNet dataset contributes significantly to semantic image understanding. Thus, by associating the progress in bridging the semantic gap with the construction of large-scale annotated datasets, we will introduce several commonly used datasets that help provide supervision to capture the semantic information.

The second aspect, optimization design under supervision, is the main focus of this paper. Although there is a vast amount of references on this topic, our goal is to shed light on the role of supervision in bridging the semantic gap. To illustrate this point, we choose two representative domains for elaboration: (1) object detection and (2) metric learning in the context of content-based image retrieval (CBIR). Both are fundamental computer vision problems. We summarize various forms of supervision used to bridge the semantic gap in the two fields, including fully-supervised learning, semi-supervised learning, weakly-supervised learning, self-supervised learning, etc.

Semantic Gap. Understanding semantics is the most fundamental step in all kinds of computer vision problems as it paves the way for general artificial intelligence. The semantic gap is a general term widely used in content-based image retrieval (CBIR). It is defined in [134] as follows. “Humans tend to use high-level concepts in everyday life. However, what current computer vision techniques can automatically extract from image are mostly low-level features. In constrained applications, such as the human face and finger print, it is possible to link the low-level features to high-level concepts (faces or finger prints). In a general setting, however, the low-level features do not have a direct link to the high-level concepts.”

Table 1: The top 1 classification accuracy for ImageNet as a function of years [121]. The ImageNet is an important dataset that drives research on object classification/recognition, and the associated image labels offer supervision to address the semantic gap problem.



Method	Top 1 ACC (%)	Top 5 ACC (%)	Params (M)	Extra data	Year
Meta Pseudo labels [123]	90.35	98.8	480	✓	2020
FixResNeXt-101 32 × 48d [151]	86.4	98.0	829	✓	2019
NASNET-A(6) [205]	82.7	96.2	88.9	✗	2017
Inception-V3 [147]	78.8	-	-	✗	2015
SPPNet [64]	72.14	91.86	-	✓	2014
Five Base + Five HiRes [67]	66.3	86.3	-	✗	2013
AlexNet [80]	63.3	84.6	60	✗	2012
SIFT + FVs [77]	50.9	73.8	-	✗	2010





Raw Media images	Descriptors Feature-vectors	Objects Prototypical combinations of descriptors	Object Labels Symbolic names of objects	Semantics Object relationships and more
	Segmented blobs, salient regions, pixel-level histograms, Fourier descriptors...			

Figure 1: Illustration of the semantic gap in multiple representation levels between raw pixels and full semantics [60]. Left to right indicates increasing levels of semantic understanding, based on information of the previous step. For example, objects are agglomeration of feature descriptors and object labels are derived based on features for the objects. The rightmost is closest to human-level understanding of the input such as object relationships and more.

The semantic gap manifests itself through different semantic understanding levels as shown in Figure 1. The raw media representation lies at the lowest level. In the context of object detection and image retrieval, the basic representation unit is the RGB pixel. At a higher level, low-level feature vectors are extracted by image analysis tools. This process is sometimes called low-level computer vision. The extracted features can be in form of segmented blobs, texture statistics, simple colour histograms, and other hand-crafted feature vectors used to represent parts or full images. As these feature descriptors are often human-engineered, they may require the domain knowledge from experts. At the next higher level, there are object representations which may be prototypical combinations of feature vectors or other more explicit representations. Once identified, objects are given symbolic labels such as object names. Labels may be general or specific, for example, an animal or a wolf. Labelling all objects does not necessarily capture the full semantics of an image since there may exist relationship between objects. Furthermore, the amount of labor required by labeling is tremendous. At the highest level as shown in Figure 1, we target at understanding the relations between objects and the holistic meaning of an image.

Semantics is a broad topic and in this survey, we choose object detection and metric learning as two examples for the following reasons:

1. We focus on resolving the gap between image pixels and semantic meanings of images, also known as the semantic gap. Although there are other vision applications such as scene graph generation [173], visual question answering [3], human-object interaction [181] that require richer semantic understanding, they are higher-level computer vision tasks that build upon object detection and metric learning. The two topics are fundamental and closer to the pixel level, where semantics is hard to capture.
2. Another reason is that there is a decent accumulation of various supervisions for the two applications, making them good choices to study the role that supervision plays in bridging the semantic gap.

That being said, there are many other useful tools such as semantic parsing [8] in the NLP domain that help understand the semantics.

Comparison with Previous Reviews. Many papers that aimed to study the semantic gap problem have been published. They are summarized in Table 2. The list includes many excellent surveys on the specific problem of image retrieval [60, 108, 22, 1], video retrieval [90], semantic segmentation [119], etc.

Recent success and dominance of deep learning based methods uphold the promise to achieve this goal. To this end, there have been many published surveys on deep learning such as the work of [7, 86, 55], and recent tutorials given at CVPR and ICCV. Although deep learning based methods have been proposed to bridge the semantic gap, we are unaware of any comprehensive survey that attempts to unify them at a higher level. In this survey, we organize papers and summarize their ideas by grouping them into different supervision forms such as fully-supervised, unsupervised, semi-supervised, self-supervised and weak-supervised etc. Another important distinction between our paper and previous ones is that we do not restrict ourselves to a specific problem but focus on resolving the semantic gap at a broader context.

Scope of our work. The central theme of our paper is supervision, which we believe is the key to semantic gap resolution. However, supervision is a broad topic and we need to limit our scope to two important problems (i.e., object detection and metric learning) as they help reveal our insights. There are still too many papers on these two topics, and compiling an extensive list of the state-of-the-art methods of both is beyond the scope of a paper of reasonable length. Instead, for domain-specific surveys, readers are referred to Tables 3 and 4, respectively.

The first selected illustrative topic is object detection. It is a fundamental computer vision problem and serves as a building block for image segmentation [59], object tracking [16, 94, 199, 175], landmark detection [187, 174], etc. Its goal is to identify the location (i.e., the coordinates of a bounding box) of an object instance and its corresponding category (for example, persons, pedestrians, cars, and animals). Previous survey papers have covered different aspects of object detection such as pedestrian detection [41, 49, 32], face detection [179, 188, 34], vehicle detection [143], gesture recognition [38, 37], text detection [184], etc. There are also a number of review papers on generic object detection methods [54, 4, 198, 102]. Here, our goal is to provide a connection between object detection and different supervisions. An example of object detection is given in Figure 2. In a fully-supervised setting, object classes and their bounding boxes are annotated in each image. However, it is expensive and often impossible to manually labor all possible objects in the real world. That is the reason other forms of supervision have to be developed. We will review various supervisions developed for object detection in Section 3.

Table 2: Summary of surveys on the semantic gap study

No.	Survey Title	References	Year	Venue	Content
1	Bridging the semantic gap in image retrieval	[200]	2002	IGI	Image retrieval
2	Bridging the semantic gap in sports video retrieval and summarization	[90]	2004	JVCI	Sports video retrieval
3	Towards bridging the semantic gap in multimedia annotation and retrieval	[158]	2006	SWAMM	Multimedia retrieval
4	Foaling the music: Bridging the semantic gap in music recommendation	[15]	2008	ISWC	Music recommendation
5	Mind the Gap: Another look at the problem of the semantic gap in image retrieval	[60]	2006	ISOP	Image-retrieval
6	Bridging the semantic gap in multimedia information retrieval: Top-down and bottom-up approaches	[61]	2015	-	Multimedia information retrieval
7	Bridging the semantic gap for texture-based image retrieval and navigation	[71]	2009	JOM	Image-retrieval
8	Bridging the semantic gap between image contents and tags	[108]	2010	IEEE MultiMedia	Image retrieval
9	ilike: Bridging the semantic gap in vertical image search by integrating text and visual features	[22]	2012	KDE	Image-text retrieval
10	A new strategy for bridging the semantic gap in image retrieval	[1]	2017	JCSE	Image-retrieval
11	Towards bridging semantic gap to improve semantic segmentation	[119]	2019	ICCV	Semantic segmentation

Table 3: Summary of surveys on object detection.

No.	Survey title	References	Year	Venue	Content
1	Survey of pedestrian detection for advanced driver assistance systems	[49]	2010	PAMI	Pedestrian detection
2	Detecting faces in images: a survey	[179]	2002	PAMI	First survey of face detection from a single image
3	A survey on face detection in the wild: past, present and future	[188]	2015	CVIU	Face detection
4	On road vehicle detection: a review	[143]	2006	PAMI	Vehicle detection
5	Text detection and recognition in imagery: a survey	[184]	2015	PAMI	Text detection
6	Object class detection: a survey	[198]	2013	ACM CS	Object detection before 2011
7	Salient object detection: a survey	[10]	2014	arXiv	A survey for salient object detection
8	A survey on deep learning in medical image analysis	[101]	2017	MIA	Object detection for medical images
9	Deep learning for generic object detection: a survey	[102]	2020	JVCI	A comprehensive survey of deep learning for generic object detection

Table 4: Summary of surveys on metric learning.

No.	Survey Title	References	Year	Venue	Content
1	Metric learning: A survey.	[81]	2012	Journal	Metric learning based on hand-crafted features
2	Distance metric learning: A comprehensive survey	[177]	2006	MSU	Traditional metric learning methods before 2006
3	Deep metric learning: A survey	[76]	2019	Symmetry	Deep learning based metric learning before 2019
4	A survey on metric learning for feature vectors and structured data	[6]	2013	arXiv	Traditional metric learning methods
5	Survey on distance metric learning and dimensionality reduction in data mining	[162]	2015	DMKD	Metric learning in data mining
6	Survey and experimental study on metric learning methods	[92]	2018	Neural networks	Experimental benchmark for metric learning
7	An overview of distance metric learning	[176]	2017	CVPR	Overview of traditional metric learning methods
8	A decade survey of content based image retrieval using deep learning	[40]	2021	CSVT	Deep learning based methods for metric learning in recent 10 years
9	A metric learning reality check	[115]	2020	ECCV	A benchmark for DML

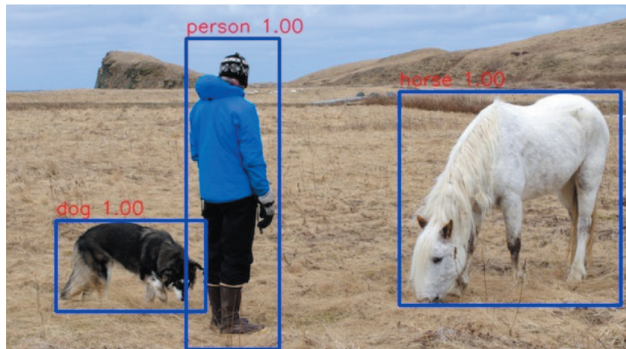


Figure 2: An object detection example.

The second exemplary topic is metric learning for image retrieval. It learns a distance metric so as to establish similarity or dissimilarity between objects and find applications in image, video and multimedia retrieval and music recommendation. While metric learning aims to reduce the distance between similar objects, it also intends to increase the distance between dissimilar objects. Typically, deep metric learning requires the class label for each individual sample. This demanding requirement prohibits its applicability in wild scenarios. Efforts have been made to relax the stringent requirement so as to accommodate other learning environments such as semi-supervised, weakly-supervised, pseudo-supervised, self-supervised and even unsupervised approaches. We will review deep metric learning methods with different supervision types in Section 4.

2 Background Review

Object Detection. Different from current deep learning based methods which extract the feature representation from images implicitly and automatically, traditional object detection methods rely heavily on hand-crafted features. A traditional object detection pipeline consists of the following three steps.

1. Extract a region of interest (say, a region that enclose objects).
2. Obtain features from the region of interest. They are often handcrafted based on the domain knowledge.
3. Classify the region of interest into a certain object class based on extracted features.

In the first step, regions of interest (ROIs) are often extracted with a sliding window approach. It requires the choice of certain hyper-parameters such as window's width, height, stride and aspect ratio. As the number of objects in the scene increases, this brute force enumeration approach can lead to a very high computational cost. Later, researchers came up with efficient yet heuristic approaches such as selective search [155], edge boxes [204], box refinement [93] to generate region proposals. The second step involves feature engineering that plays a crucial role in the performance. One seminal work is the Scale Invariant Feature Transform (SIFT) [107]. It was designed to be robust against changes in translation, scale, rotation, illumination, viewpoint and occlusion. Other local representative descriptors [111] include Haar-like features [159], local binary patterns (LBP) [118] and region covariances [153]. The histogram of oriented gradients (HOGs) [27] is an important improvement over SIFT and offers a better object descriptor. The HOG feature is robust against local deformation and illumination, and it has been widely used in classical object detectors. The last step is classification based on features of each ROI. Most commonly used classifiers include the support vector machine (SVM) [25], AdaBoost [45] and random forest (RF) [144].

One famous example of the three-step pipeline is the Viola-Jones face detector [159]. It adopts a sliding window approach to check if a face object is included in the window. To improve the detection speed, it uses an Adaboost training approach and cascades classifiers to improve the detection performance. The deformable part model (DPM) [44] offers another milestone in the traditional object detection framework. It consists of a root-filter and multiple part-filters. DPM improves HOG using hard negative mining, bounding box regression and context priming. It was the champion solution in the Pascal-VOC Challenge from 2007 to 2009 [43]. The cascaded pipelines of "hand-crafted feature description" followed by "discriminative classification" dominated many computer vision tasks, including object detection, for years. Even with significant advancement, there is a substantial gap between the classical object detector and human recognition capability. This gap is attributed to two main barriers: the limited representation capability of hand-crafted features and lack of sufficient supervision. Deep learning can learn powerful features automatically to overcome the first barrier. The construction of larger and larger labeled datasets addresses the second barrier.

Metric Learning. Metric learning is a branch of machine learning. It learns a distance metric that establishes similarity and dissimilarity between objects from training images. The objective is to reduce the distance between similar objects while increasing that between dissimilar ones. The task is also known as similarity learning and it is most commonly used in image retrieval [87, 178, 66, 110, 68], person re-identification [172, 185, 97] and face verification [116], etc.

For a given query image, a content-based image retrieval (CBIR) system [134] return a ranked list of images from the database based on a similarity measure between the query and retrieved images. CBIR is a challenging problem since it is often that many (or even all) of those returned images do not look similar to the query one from a human perspective. This is because that most similarity metrics refer to distances of low-level features. They do not correlate well with semantic similarity perceived by humans.

Traditional CBIR research focused on two areas: feature design and distance metric selection. Research on the application of hand-crafted features to CBIR was rich. SIFT [107] and LBP [118] were widely used features. A histogram-based similarity measure was proposed in [145] for image retrieval. The K-means clustering approach was used in [137] to discover the patterns of data in low-level feature space using the color information. A nonlinear mapping approach based on sparse kernel learning was studied in [113]. Other features were designed based on the prior knowledge and domain expertise for specific application. For example, LOMO [97, 35] was developed to deal with illumination and viewpoint changes for the matching of person images. As to distance metrics, common choices include Euclidean, Mahalanobis [23, 28, 170], and Kullback-Leibler distances. Higher performance may be achieved by mapping the problem to a non-linear space through kernel methods. These non-linear methods are often used in combinationn of regularization techniques [72] to avoid overfitting.

Large-Scale Labeled Datasets. Resolution of the semantic gap has been a long standing problem. There is a substantial progress on this topic in the last decade. This progress is attributed to the realization of the importance of supervised learning.

While the semantic gap is bridged by supervision, the key to supervised learning is the availability of large-scale human annotated datasets. Judged by today’s standard, the sizes of labeled datasets were quite small before 2010. This practice can be attributed to several reasons. First, it is about the labeling cost. The labor required to annotate collected data is substantial. Second, most traditional methods do not scale well with the data size. When most solutions do not work well for small datasets, the motivation in building larger datasets would not be strong. Third, since people can understand semantic meanings from a small set of examples (i.e., weakly-supervised learning), it is natural to expect powerful vision algorithms to do the same. For all these reasons, the importance of “large-scale supervision” was not appreciated and practiced until the last decade.

The situation began to change with the introduction of the ImageNet dataset [29], which was viewed as the engine to drive deep learning in early 10s. That is, deep learning has gained widespread attention and popularity after [80] achieved record-breaking image classification accuracy in the Large Scale

Visual Recognition Challenge (ILSVRC) [135] in 2012. Although deep learning provides a mechanism in extracting powerful representative features, feature extraction is not the main objective of deep learning but a byproduct. Deep learning relies heavily on supervision. It attempts to build a mapping from images to labels by certain neural networks. In other words, it uses human labels as the ground truth and provides a nonlinear mechanism that minimizes the error between the predicted and true labels.

The chase of more and more annotated data in today’s machine learning community is a clear evidence of supervision’s role in bridging the semantic gap. A tremendous amount of efforts have been spent in data collection and labeling nowadays. In the following, we will highlight several datasets that are critical to the development of object detection and metric learning techniques.

Four datasets are commonly used for generic object detection: PASCAL VOC ([43], ImageNet [29]), MS COCO [100], and Open Images [78]. The attributes of these datasets are summarized in Table 5. The selected samples are shown in Figure 3. Several criteria are used in evaluating the performance of an object detector, including precision, recall, model sizes, and inference speed measured by frames per second (FPS). While the average precision (AP) that combines precision and recall is used to evaluate the performance for a specific category, the mean AP (mAP) averaged over all categories is used as the measure of performance over all categories. For more details, readers are referred to [100].

There are also four datasets that are commonly been used in metric learning. They are: CUB-200-2011 [160], CAR-196 [79], Stanford Online Shopping [117] and Market-1501 [201]. The first two focus on fine-grained object category retrieval. The last two are instance-level retrieval datasets. Market-1501 is one of the largest person-reidentification benchmark dataset. Detailed statistics are shown in Table 6 and exemplary images are shown in Figure 4. Precision@ k , denoted by $P@k$, is a popular metric in metric learning. It indicates the number of relevant images among the top k retrieved images. If there are R images that belong to the same class as the query, the R-precision (RP) measures the percentage of correct retrievals among the top R retrieved results. Another recently proposed metric is MAP@R [115], that combines the idea of mean average precision with RP to offer a more accurate performance measure.

3 Supervision for Object Detection

3.1 Full Supervision

Deep learning methods have been extensively developed for the fully-supervised object detection task [53, 64, 51, 129, 136]. Research on this topic has reached quite a mature stage. Generally speaking, deep-learning-based object detection

Table 5: Commonly used databases for object detection.

Dataset name	Total images	Categories	Images per category	Objects per image	Image size	Started year	Highlights
PASCAL VOC (2012)	11,540	20	303 ~ 4087	2.4	470 × 380	2005	Covers only 20 categories that are common in everyday life; Large number of training images; Close to real-world applications; Significantly larger intraclass variations; Objects in scene context; Multiple objects in one image; Contains many difficult samples.
ImageNet	14M	21,841	—	1.5	500 × 400	2009	Large number of object categories; More instances and more categories of objects per image; More challenging than PASCAL VOC; Backbone of the ILSVRC challenge; Images are object-centric.
MS COCO	328,000+	91	—	7.3	640 × 480	2014	Even closer to real world scenarios; Each image contains more instances of objects and richer object annotation information; Contains object segmentation notation data that is not available in the ImageNet dataset.
Places	10M	434	—	—	256 × 256	2014	The largest labeled dataset for scene recognition; Four subsets Places365 Standard, Places365 Challenge, Places 205 and Places88 as benchmarks.
Open Images	9M	6000+	—	8.3	varied	2017	Annotated with image level labels, object bounding boxes and visual relationships; Open Images V5 supports large scale object detection, object instance segmentation and visual relationship detection.

Table 6: Commonly used databases for metric learning.

Dataset name	Total images	Categories	Training images	Testing images	Retrieval level	Year	Highlights
CUB-200-2011	11,788	200	5,864	5,924	Object	2011	Fine-grained bird retrieval dataset
CAR-196	16,185	196	8,054	8,131	Object	2013	Fine-grained car retrieval dataset
Stanford online	120,053	22,634	59,551	60,502	Instance	2016	Covers a variety of online shopping instances.
Market-1501	32,000	1501	750 id	751 id	Instance	2015	One of the most widely used person re-identification dataset



Figure 3: Selected sample images for popular object detection datasets [43, 29, 100, 78].

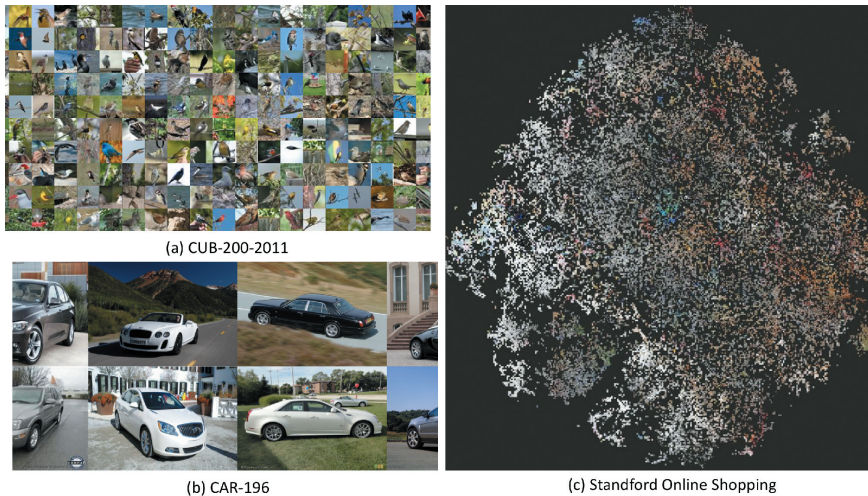


Figure 4: Selected sample images for popular metric learning datasets [160, 79, 117, 201]. Subfigure (c) is a tSNE [157] visualization of the dataset [117].

methods can be categorized into two types: two-stage detection and one-stage detection. Recently, there’s an emerging line of transformer-based works [13, 203, 26] which approach object detection as a direct set prediction problem. We elaborate representatives for each of the category below.

Two-Stage Detection. Two-stage detection methods consist of two stages in cascade: (1) the region proposal stage and (2) the object classification stage. The common pipeline includes: category-independent region proposals¹ are generated from an image, CNN features are extracted from these regions, and then category-specific classifiers are used to determine the category label of each proposal.

¹Object proposals, also called region proposals or detection proposals, are a set of candidate regions or bounding boxes in an image that may potentially contain an object.

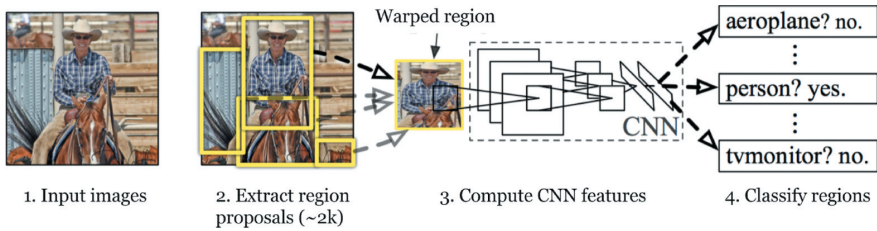


Figure 5: The pipeline of a two-stage object detection framework [53].

Inspired by the impressive image classification performance achieved by the AlexNet [80] and the success of selective search in finding region proposals with hand-crafted features ([155, 53, 52]) were among the first to explore CNNs for generic object detection and developed RCNN as shown in Figure 5. The training of an RCNN framework consists of the following steps:

1. **Region proposal selection** Class agnostic region proposals, which are candidate regions that might contain objects, are obtained via selective search.
2. **Region proposal processing** Selected region proposals are cropped from the image and warped into the same size. They are used as the input to finetune a CNN model pre-trained by a large-scale dataset such as ImageNet. In this step, a region proposal with its IOU against a ground truth box greater than 0.5^2 is defined as a positive for the ground truth class and the rest as negatives.
3. **Class-specific SVM classifiers training** A set of class-specific linear SVM classifiers are trained using the fixed length features extracted by the CNN, replacing the softmax classifier learned by finetuning. For the training of SVM classifiers, positive examples are the ground truth boxes for each class. A region proposal that has less than 0.3 IOU with all ground truth instances of a class is negative for that class. Note that the positive and negative examples used for training SVM classifiers are different from those for finetuning the CNN.
4. **Class-specific bounding box regressor training** Bounding box regression is learned for each object class with CNN features.

There are variants of RCNN for better performance. Two noticeable ones are Fast RCNN [51] and Faster RCNN [129] as shown in Figure 6. Fast RCNN improves both detection speed and accuracy of RCNN. Rather than separately training a softmax classifier, SVMs, and bounding box regressors as done

²This is a commonly used practice such as in MSCOCO dataset [100].

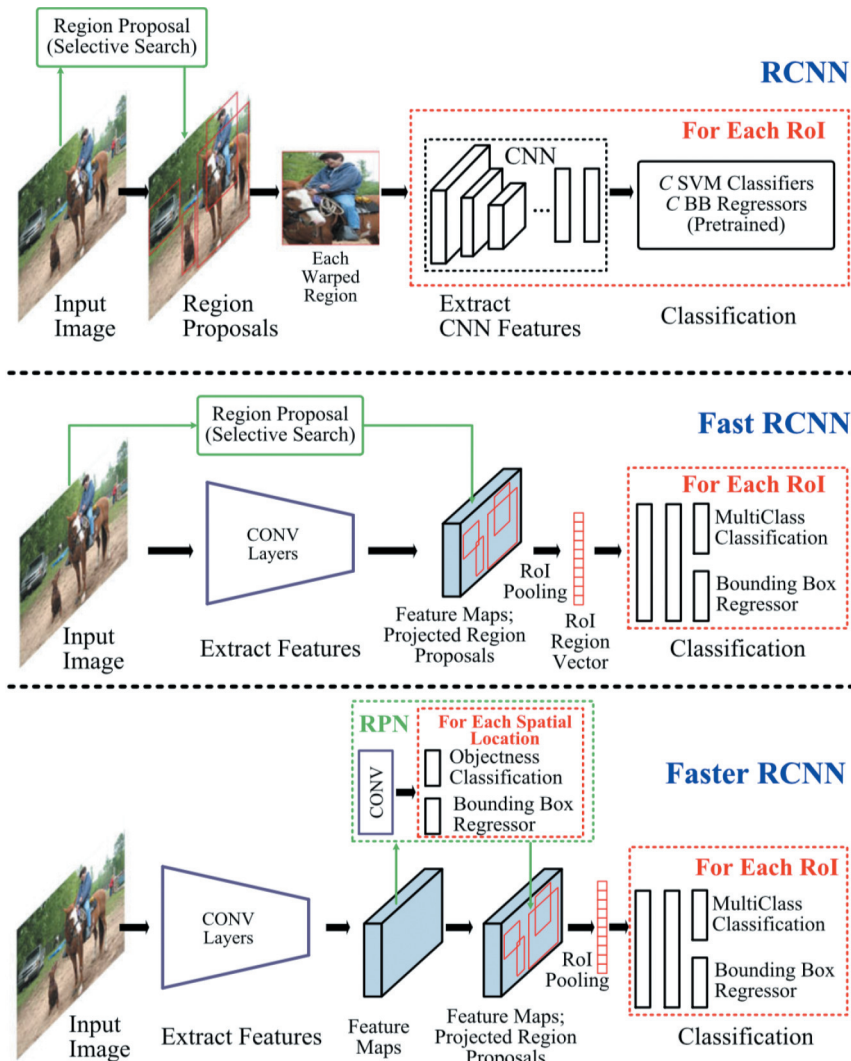


Figure 6: The system diagrams of three two-stage object detection methods [102]: RCNN (top), Fast RCNN (middle), and Faster RCNN (bottom).

in RCNN, Fast RCNN enables end-to-end detector training by developing a streamlined training process that simultaneously learns a softmax classifier and class-specific bounding box regression. The core idea of Fast RCNN is to share the feature extraction process among different region proposals. Fast RCNN improves efficiency of RCNN considerably, typically 3 times faster in training and 10 times faster in testing, and there is no storage required for feature

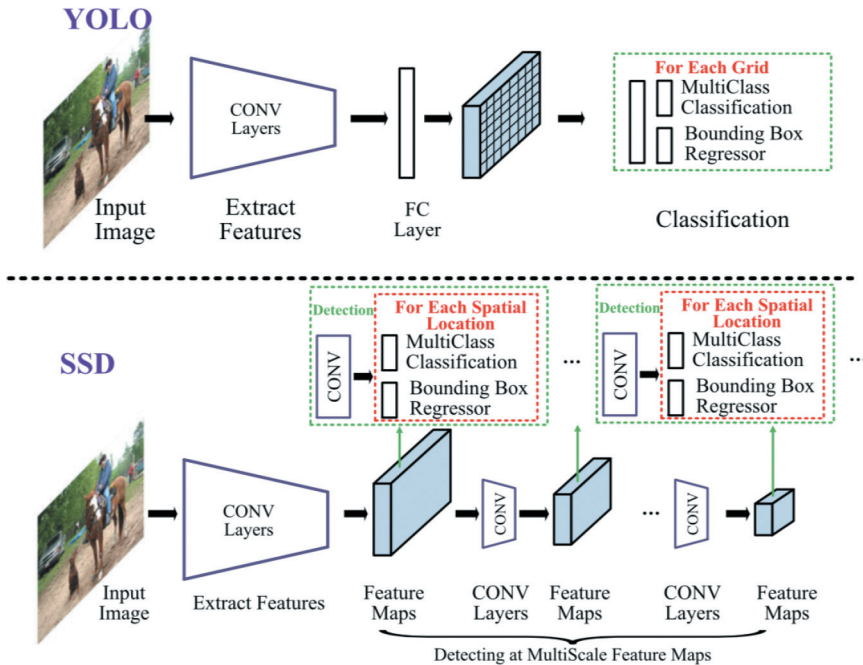


Figure 7: The system diagrams of two one-stage object detection networks [102]: YOLO (top) and SSD (bottom).

caching. Faster RCNN offers an efficient and accurate region proposal network (RPN) in generating region proposals. It utilizes the same backbone network but exploits features from the last shared convolutional layer to accomplish the task of RPN for region proposal generation and the task of Fast RCNN for region classification.

The two-stage region-based pipeline offers state-of-the-art object detection performance as evidenced by the fact that leading results on popular benchmark datasets are all based on Faster RCNN [129]. Nevertheless, region-based methods are computationally costly for mobile/wearable devices with limited storage and computational resources. Instead of optimizing individual components of the complex region-based pipeline, researchers looked for an alternative that detects objects directly without the region proposal step.

One-Stage Detection. By one-stage detection, we refer to an architecture that predicts class probabilities and bounding box sizes and locations from full images with a single feed-forward CNN in a monolithic setting. It can be optimized end-to-end directly on detection performance. DetectorNet [146] was among the first in exploring this new direction. However, since the network

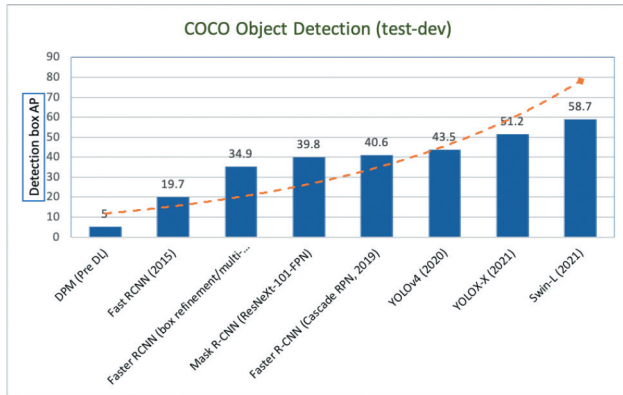


Figure 8: Progress of object detection performance on Microsoft COCO over years, where results are quoted from ([120]).

needs to be trained per object type and mask type, it does not scale well as the number of classes increases. [128] proposed YOLO (You Only Look Once), which is a unified detector casting object detection as a regression problem from image pixels to spatially separated bounding boxes and associated class probabilities. Unlike the two-stage detection, YOLO predicts detections based on features from local regions of multiple sizes. Specifically, YOLO divides an image into an $S \times S$ grid, each predicting C class probabilities, B bounding box locations, and confidence scores. By eliminating the region proposal generation step entirely, YOLO is fast by design, running in real time at 45 FPS. Fast YOLO [138] can even reach 155 FPS. To preserve real-time speed without sacrificing much detection accuracy, [103] proposed SSD (Single Shot Detector), which is faster than YOLO and with an accuracy competitive with region-based detectors such as Faster RCNN [129]. SSD effectively combines ideas from RPN in Faster RCNN and YOLO multiscale CONV features to achieve fast detection speed, while still retaining high detection quality. Like YOLO, SSD predicts a fixed number of bounding boxes and scores, followed by a non-maximum-suppression (NMS) step to produce the final detection. The system diagrams of YOLO and SSD are shown in Figure 7 for comparison.

Recently, there’s a growing trend in applying transformers to the computer vision tasks [13, 33, 150], among which DETR [13] is a representative for object detection. Instead of generating “proposals”, it patchifies the given image as tokens and feeds them to the vision transformer. While previous detectors rely on NMS as a post-processing step, DETR casts object detection as a set prediction problem and leverages the Hungarian algorithm to match the box predictions with ground-truth boxes during training. DETR simplifies traditional object detection pipeline and obtains 42 AP on COCO using a Resnet50 backbone.

Further Performance Improvement. We show the detection performance improvement over years with respect to the Microsoft COCO challenge in Figure 8. As to the object detection task in the open image challenge, the current leader [106] achieved 58.7 box AP in the public leader board. It proposes a hierarchical transformer whose representation is computed with Shifted windows. Representative methods benchmarked include Fast RCNN [51], Faster-RCNN [129], FPN [50], Deformable Faster RCNN [130], Cascade RCNN [11], Mask-RCNN [63] and YOLO families [128]. Generally speaking, the backbone network, the detection pipeline and the availability of large-scale training datasets are three most important factors in further detection accuracy improvement. Besides, ensembles of multiple models, the incorporation of context features, and data augmentation all help achieve better accuracy.

3.2 Weak Supervision

Object detectors are trained without bounding box annotations in weak supervision detection (WSD), where only image-level labels are used. The main challenge of WSD is object localization since a label may refer to any object in the image. This problem is typically addressed using multiple instance learning, which is a well-studied topic [9, 24, 161]. Although image-level labels are easier to collect than bounding boxes, they still require manual efforts. Besides, they are often limited to a pre-defined taxonomy.

Some recent work adopts captions, which are often freely available on the web. Learning object detection from captions has been studied but at a limited scale. CAP2Det [182] parses captions into a multi-label classification targets and, then, these labels are used to train a WSD model. However, it requires image-level labels to train the caption parsers. Besides, it is under the constraint of a closed vocabulary. Another WSD model was trained in [2] based on a predefined set of words in captions. It is similar to a closed vocabulary, yet the rich semantic content of captions is discarded. In contrast, research in [141] and [183] aims to discover an open set of object classes from image-caption corpora, and learns detectors for each discovered class.

One shortcoming of WSD methods is their poorer object localization accuracy. Object recognition and localization are disentangled into two independent problems in [189]. It learns object localization using a fully annotated dataset from a small subset of classes and conducts object detection using open-vocabulary captions.

3.3 Mixed Full/Weak Supervision

Mixed supervision has been studied to exploit both full supervision and weak supervision. Most mixed supervision methods need bounding box annotations for all main classes and apply weak supervision to auxiliary classes [168, 47, 127].

For example, by following the transfer learning framework, one can transfer a detector trained on supervised base classes to weakly supervised target classes [65, 149, 154]. These methods usually lose performance on target classes.

One common limitation of mixed-supervised methods is that they require image-level annotations within a predefined taxonomy so that they learn the predefined classes only. To address it, one recent work [31] exploits supervision from captions that are open-vocabulary and more prevalent on the web. Instead of training on base classes and transferring to target classes, it uses captions to learn an open-vocabulary semantic space that includes target classes, and transfers that to the object detection task via supervised learning.

3.4 Zero-shot Detection

Zero-shot object detection (ZSD) aims to generalize from seen object classes to unseen ones by exploiting zero-shot learning techniques (e.g., word embedding projection [46]) for object proposal classification. There exist however different views on ZSD. According to [5], the main challenge of ZSD lies in modelling the background class since it is difficult to separate from unseen classes. The background was treated as a mixture model in [5]. It was furthermore improved by introducing the polarity loss [125]. On the other hand, it was argued in [202] that the key challenge of ZSD is the generalization capability of object proposal models. To tackle with it, they employed a generative model to hallucinate unseen classes and augment seen examples in the proposal model training process.

Nevertheless, there is still a significant gap in the performance due to the unnecessarily harsh constraint; namely, not having any example of unseen objects, and having to guess how they look like solely based on their word embeddings [5, 125] or textual descriptions [96]. This has motivated researchers to simplify the task by making more assumptions such as the availability of test data during training [126] or the availability of unseen class annotations to filter images with unseen object instances [57]. Since datasets with natural, weak supervision are abundant and cheap, an alternative was proposed to utilize image-caption datasets in [189], which covers a larger variety of objects with an open vocabulary.

4 Supervision for Metric Learning

4.1 Full Supervision

Metric learning attempts to map image data to an embedding space, where images of similar semantic content are closer while those of dissimilar semantic meaning are farther apart. The embedding learned in this way captures

semantics intuitive to human understanding which are initially not obvious in the pixel form. In general, this objective can be achieved by leveraging embedding and/or classification losses.

Embedding Losses. The embedding loss operates on the relationship between samples in a batch while the classification losses include a weight matrix that transforms the embedding space into a vector of class logits. Typically, embeddings are preferred when the task is in form of information retrieval whose goal is to return a data sample that is most similar to the query one. A specific example is image retrieval, where the input is a query image and the output is the most similar image in a database. Another application context is open-set classification where the test set and the training set classes are disjoint, and there are cases no proper classification loss can be defined. For example, when constructing a dataset, it might be difficult (or costly) to assign the class label to each sample. It might be easier to specify the relative similarities between samples in form of pair- or triplet-relationship. Pairs and triplets can provide additional information for existing datasets. Since both do not have explicit labels, embedding losses become a suitable choice. Pair and triplet losses provide the foundation to two fundamental metric learning approaches (See Figure 9).

Contrastive Loss. A classic pair-based loss is the contrastive loss [58] in form of

$$L_{\text{contrastive}} = [d_p - m_{\text{pos}}]_+ + [m_{\text{neg}} - d_n]_+,$$

where notation $[x]_+$ denotes $\max(x, 0)$. In the above equation, it attempts to make the distance between positive pairs d_p smaller than threshold m_{pos} , and the distance between negative pairs d_n larger than threshold m_{neg} . A theoretical downside of this method is that the same distance threshold is applied to all pairs even though there may be a large variance in their similarities and dissimilarities. The triplet margin loss [170] is developed to address this issue.

Triplet Loss. A triplet consists of an anchor input, A , a positive input, P , and a negative input N , where the anchor is more similar to the positive than the negative. The triplet margin loss is used to ensure that the anchor-positive distance (d_{ap}) is smaller than the anchor-negative distance (d_{an}) by a predefined margin (m). The triplet loss function can be written in form of

$$L_{\text{triplet}} = [||f(A) - f(P)||^2 - ||f(A) - f(N)||^2 + m]_+,$$

where f is an embedding function. This triplet loss places fewer restrictions than the contrast loss in the embedding space. It allows a learned model to account for the variance in interclass dissimilarities.

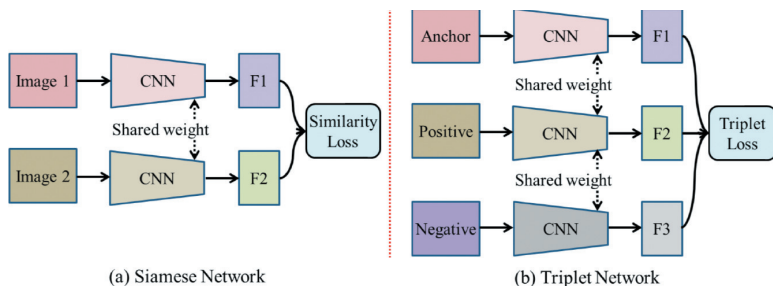


Figure 9: Comparison of the similarity loss and the triplet loss using the siamese network [40]. For Siamese Network, it optimizes to increase the similarity between positive pairs and decrease the similarity between negative pairs. The Triplet Network enforces the distance between the anchor and the positive to be smaller than that between the anchor and negative.

Other Loss Functions. A wide variety of loss functions has been defined based on these fundamental concepts. For example, the angular loss [167] is a triplet loss where the margin is based on the angles formed by the triplet vectors. The margin loss [171] modifies the contrastive loss by setting

$$m_{pos} = \beta - \alpha, \text{ and } m_{neg} = \beta + \alpha,$$

where α is fixed, and β is learnable. Other pair losses are based on the softmax function and LogSumExp, which is a smooth approximation of the maximum function. Specifically, the lifted structure loss [117] is the contrastive loss but with Log-SumExp applied to all negative pairs. The N-Pairs loss [140] applies the softmax function to each positive pair relative to all other pairs. It is also known as InfoNCE [156] and NT-Xent [18]. The tuplet margin loss [186] combines Log-SumExp with an implicit pair weighting method while the circle loss [142] weighs each pair’s similarity by its deviation from a pre-determined optimal similarity value. A general weighting framework was presented in [169] to understand recent pair-based loss functions. In contrast with pair and triplet losses, FastAP [12] attempts to optimize for average precision within each batch using a soft histogram binning technique.

Classification Losses. Classification losses are obtained by including of a weight matrix, where each column corresponds to a particular class. In most cases, the training process consists of multiplying weight matrix with embedding vectors to obtain logits, and then applying a certain loss function to the logits. The most straightforward one is the normalized softmax loss [164, 104, 190]. It is identical with the cross entropy loss with L2-normalized columns of the weight matrix.

One variant is ProxyNCA [114], where the cross entropy loss is applied to the Euclidean distances, rather than the cosine similarities, between embed-

dings and the weight matrix. A number of face verification losses modified the cross entropy loss with angular margins in the softmax expression. For example, SphereFace [104], CosFace [163, 166] and ArcFace [30] apply multiplicative-angular, additive-cosine and additive-angular margins, respectively. It is interesting to note that many metric learning papers leave out face verification losses from their experiments although they are not face-specific. The Soft-Triple loss [124] takes a different approach by expanding the weight matrix to have multiple columns per class. It has more flexibility in modeling class variances.

4.2 Other Forms of Supervision

Embedding of semantic information is hard to learn when the amount of labeled data is limited or when the data is imbalanced, which is often the case in real-world scenario. The research community tries to address the semantic issues by either adopting the weakly labeled data, partially labeled data or unlabeled data, leading to different supervisions discussed below.

Weak Supervision. Weakly supervised approaches have been explored for the image retrieval task [148, 56, 48, 95]. For example, Tang et al. [148] proposed a weakly-supervised multimodal hashing method that exploits local discriminative and geometric structures in the visual space. [56] performed pre-training in weak supervision mode and finetuned the network in supervision mode. [48] developed a weakly supervised deep hashing method that used tag embeddings for image retrieval with the word2vec semantic embeddings. [95] developed a semantic guided hashing network for image retrieval by employing the weakly-supervised tag information and inherent data relations simultaneously.

Semi-Supervision. The semi-supervised approaches generally use a combination of labeled and unlabelled data in feature learning. A semi-supervised deep hashing framework was proposed for image retrieval from labeled and unlabeled data in [192]. It uses labeled data for empirical error minimization and both labeled and unlabeled data for embedding error minimization. The generative adversarial learning approach was also utilized in semi-supervised deep image retrieval [165, 73]. A teacher-student semi-supervised image retrieval method was presented in [197], where the pairwise information learned by the teacher network is used as the guidance to train the student network. Pseudo labels are another source of supervision falling into semi-supervised regime. For example, [69] generates pseudo labels based on the pretrained VGG16 features via k-means clustering. In [36] a self-training framework, SLADE, is proposed to improve retrieval performance by leveraging additional unlabeled data. It first train a teacher model on

the labeled data and use it to generate pseudo labels for the unlabeled data. It then train a student model on both labels and pseudo labels to generate final feature embeddings. The framework significantly improves existing state-of-the-art.

Self Supervision. There has been an arising amount of interest in self-supervised learning [62, 18]. It is similar to unsupervised learning but apply self-generated pseudo-labels to the data during the training process. Self supervision is often achieved by clever usage of data augmentation or information from other modalities. For example, [112] defines a set of pretext tasks for learning invariant feature representation. In [14], an online vision transformer was asked to predict the output of a target vision transformer, whose input is an augmentation of the first transformer’s input. As this requires no annotations, it is self-supervised and exhibits superior performance when applied to image retrieval. Self-supervised learning also proves useful for initializing deep metric learning embedding [36], video retrieval [191], and cross-image retrieval [91].

No Supervision. Though supervised models have shown promising performance in image retrieval, it is always difficult to get labeled large-scale data. Thus, unsupervised models have been investigated and they do not require class labels to learn features. Generally, unsupervised models enforce the constraints on hash codes and/or generate the output to learn features. [42] used deep networks in an unsupervised manner to learn hash codes with the help of constraints such as the quantization loss, balanced bits and independent bits. [70] utilized deep networks coupled with unsupervised discriminative clustering to learn the description in an unsupervised manner. [122] used an unsupervised convolutional kernel network to learn convolutional features for image retrieval. They applied it to patch retrieval as well.

[98] imposed constraints (for example, the minimal quantization loss, evenly distributed codes, and uncorrelated bits) on an unsupervised deep network and proposed a solution, called DeepBit, for image retrieval, image matching and object recognition applications. DeepBit has a two-stage training process. In the first stage, the model is trained with respect to above-mentioned objectives. To improve its robustness, the network is finetuned in the second stage based on rotation data augmentation. The analysis of DeepBit is given in [99]. However, DeepBit suffers from severe quantization loss due to rigid binarization of data using the sign function without considering its distribution property. To tackle the quantization problem of DeepBit, a deep binary descriptor with multiquantization was proposed by [39]. It is achieved by jointly learning the parameters and the binarization functions using a K-AutoEncoders (KAEs) network.

5 Future Research Directions

Supervision has been dominated in form of data labeling in the last decade. However, this form appears to be quite limited. Humans learn semantics and knowledge from a wide range of resources, for example, domain knowledge priors, correlation from different domains and modalities, etc. Although it is still a mystery how humans learn semantic meanings from the real world, it is anticipated that supervision on machines will appear in richer form. In this section, we present two general directions for future research.

5.1 *Interpretable and Modularized Learning*

Interpretability and modular design are pillars to the construction, debugging and maintenance of next-generation artificial intelligence (AI) systems. Although deep learning is the dominant methodology in providing the mapping between image pixels and semantics nowadays, it is neither interpretable nor modularized and we anticipate the same mapping to be achieved by other alternatives.

One emerging alternative is successive subspace learning [83, 82, 21, 84, 85, 131]. Simply speaking, SSL is a light-weight unsupervised data embedding (or feature learning) method and it can be applied to different data types (e.g., images, point-clouds, voxels, etc.) The SSL pipeline consists of a sequence of joint spatial-spectral transforms in cascade with PCA-like transform kernels. They are rigorously derived using statistical properties of data units such as pixels, voxels and points. SSL-based embedding is data driven and repeatable. The SSL pipeline can be connected to a classifier (for example, the random forest, the support vector machine or the extreme gradient boosting classifier, etc.) or a regressor (for example, the linear regressor, the logistic regressor, the support vector regressor, etc.) for final decision.

The representations associated with SSL are unsupervised, interpretable, modularized, robust to perturbations, effective (i.e., a small embedded dimension) and efficient (a smaller model size and low embedding complexity). Since end-to-end optimization is completely abandoned in SSL, its training complexity is significantly lower. It can be implemented on low-cost CPUs. The sizes of SSL models are significantly smaller than those of DL-based models; thus, suitable for mobile and edge computing. From the angle of supervision, SSL can incorporate priors conveniently, and fine-tune the AI system with new observations on the fly.

SSL-based solutions find applications in object classification [19, 20, 109, 180], fake face image detection [17], face gender classification [132], low-resolution face recognition [133], joint compression and classification [152], point cloud classification and registration [75, 74, 194, 196, 195], image and

texture generation [88, 89], anomaly detection [193] and medical image classification [105].

5.2 Intelligence Gaps

Intelligence gaps is a collection of three characterized aspects to target, in order to reach human-level semantic understanding from raw data perception. The three aspects envision three increasing levels of semantic understanding which we examine how to fill in below.

Gap between Signals and Semantic Units. Humans have sensors such as eyes, ears, nose, skin, etc. to receive signals (or stimuli) from the external world. They include visual, audio, smell, pressure, temperature, etc. Signals need to be converted to compact representations for future processing in machines - known as “embedding”. We have witnessed rapid progress in signal/data embedding. There are a few criteria in evaluating embedding schemes: interpretability, supervision degree, sensitivity, effectiveness and efficiency. Most today’s embedding methods rely on deep learning. They are far from ideal according to these criteria. A new signal embedding idea is to exploit statistical properties of data units (for example, pixels, vertices, and points) in an unsupervised feedforward fashion based on SSL as discussed in Section 5.1. It is interpretable, robust to perturbations, effective (i.e., smaller embedded dimension), efficient (smaller model size and lower complexity), and suitable for multi-tasking.

Furthermore, two challenges in signal embedding worth further study: (1) attention and (2) multi-modal data representation. Both machines and human brains have limitations on processing speed, memory and communication capacity. Attention is needed to enable an intelligent system to process the most relevant input within its limits. Attention is often derived by end-to-end optimization nowadays (e.g., visual saliency in computer vision and transformers in natural language processing). Yet, attention can be easily fooled with small perturbation. For example, it can be shifted from one region to another in an image by manipulating a few pixels - leading to a totally different outcome. Adversarial attacks impose a major threat in real-world applications. Interpretable and robust attention is essential in next generation AI, which will be assisted by semantic scene and object segmentation. For multi-modal data representation, subspace decomposition may be leveraged. That is, we may represent audio, image, video, and 3D data in their individual subspaces and select a suitable combination and use the direct sum of these subspaces to construct a multi-modal space. Each subspace can be updated independently and combined efficiently and dynamically in response to different needs.

Gap between Semantic Units and Knowledge. Semantic units are segmented for ease of re-composition. The study of their relationship yields a richer information space. For example, the WordNet contains relations between numerous words so as to result in a huge graph. Knowledge presents the highest abstraction level of human cognition. Besides knowledge representation and acquisition, humans can infer missing information and discover knowledge that are not directly available.

Generally, one can construct individual knowledge graphs based on existing databases in various domains and then combine them into larger heterogeneous graphs. Knowledge graphs will be a central piece of the next generation AI. There are open problems to be addressed, including scalability, ambiguity resolution, semantic matching, path finding/completion, new entity discovery, hidden relation extraction, dynamic graph evolution, etc.

To tackle with scalability, a decomposition and re-composition methodology through interaction of semantic and knowledge spaces could be a direction to explore. For example, today's CNNs recognize cars of different colors and models from various angles through numerous labeled car images. Yet, this is not how humans acquire the knowledge of "cars". Humans decompose cars into semantic units such as body, wheels, doors, windows, lights, etc. and use them to form the knowledge of "cars". The decomposition/re-composition process enables humans to learn cars with fewer examples. The success lies in the interaction of the 3D car structure (knowledge) and the projected 2D car images (semantic units). Also, it is challenging to recognize small components of cars such as wheels, doors, windows, lights, etc. alone. Yet, the 3D structural knowledge of cars can help trace/confirm the parts and make their recognition easier. The same principle applies to human perception on objects with occlusion. We can recognize occluded objects if occlusion is not severe. Generally, one recognizes objects through their salient regions and, then, their parts through the assistance of knowledge graphs.

Ambiguity resolution can be done using the context information in the knowledge space. Ontology is essential to human knowledge acquisition, organization, and learning. Hierarchical categorization is more stable and easier to update. Today's knowledge graphs are flat without ontology incorporated. Lack of knowledge hierarchy makes the representation difficult to scale up. It may be feasible to enforce the ontological relationship in sub-space decomposition. That is, a high-dimensional knowledge space will be decomposed into a direct sum of multiple low-dimension knowledge subspaces. The core knowledge, which is stored in a low-dimension subspace, should be more stable and error resilient with less frequent update. The refined knowledge in specific domains is stored in other low-dimensional subspaces. They will be updated more frequently and optimized locally. Unequal security can be applied to protect different subspaces from attacks depending

on their importance. Knowledge space decomposition and re-composition provides flexibility in face of a rapidly changing, stochastic and adversarial environment. Mathematically, this decomposition can be achieved by tensor operations.

Gap between Knowledge and Concept/Decision. Knowledge is what we know. It's the accumulation of past experience and insight that shapes the lens by which we interpret, and assign meaning to, information. In psychology, decision-making is regarded as the cognitive process resulting in the selection of a belief or a course of action among several possible alternative options. Logical inference is the basis of human reasoning. Although this is analogous to path finding in knowledge graph, path finding is often used to find the relationship between two entities rather than two concepts. Mathematical proof based on computer enumeration exists. Yet, it does not have the ability to infer from one concept to another.

To narrow the gap, we may construct the concept (or rule) graph whose nodes represent different concepts. For example, from two concepts "a car runs faster than a horse" and "a horse runs faster than a man", we infer that "a car runs faster than a man" through the transitive law. A concept graph in a general domain could be too complex to build. Yet, it could be feasible to do it in a special domain. For example, if we focus on "I for health care", the number of concepts is much smaller. There is difference between the concept graph and the traditional expert system. An expert system does not have links between nodes while links in the concept graph introduce the logical relationship between two concepts. Common sense reasoning and rules discovery/creation are possible through inductive learning, that is, path finding in domain-specific concept graphs.

Humans react to external stimuli with responses such as action, decision and planning. Rational responses are knowledge-based. Action/decision is often related to penalties and/or rewards. Reinforcement learning is developed with this principle via cost function definition and optimization. This proves to be effective in gaming (for example, chess and go). Game theory offers an alternative optimal decision process among independent and competing actors in strategic settings.

Generalization of reinforcement learning and game theory to real-world situations is however non-trivial since it is difficult to define proper cost functions. Furthermore, human behavior involves intuition, instinct, psychological factors and constraints (e.g., faith and ethics), which are difficult to model. For the next generation AI to be fully autonomous, we need a clearly defined goal. For example, medical diagnosis can be conducted by AI automatically while medical treatment will be determined by AI and humans jointly since the latter involves human factors. Also, it is the human who should take the ultimate responsibility.

6 Conclusion

Resolving the semantic gap is a topic that has attracted growing attention in the artificial intelligence community. Unlike previous papers, this survey drew experiences from two fundamental computer vision problems: object detection and metric learning in image retrieval. The central theme was on the role of supervision, which was accomplished by “data annotation schemes” and “design of loss functions”. We organized the survey by various supervision forms. Furthermore, we offer a broader perspective on intelligence gaps and discuss a couple of ideas in resolving these gaps to shed light on future research directions.

References

- [1] M. A. Alzubaidi, “A New Strategy for Bridging the Semantic Gap in Image Retrieval,” *International Journal of Computational Science and Engineering*, 14(1), 27–43.
- [2] E. Amrani, R. Ben-Ari, T. Hakim, and A. Bronstein, “Learning to Detect and Retrieve Objects From Unlabeled Videos,” in *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, IEEE, 2019, 3713–7.
- [3] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, “Bottom-up and Top-Down Attention for Image Captioning and Visual Question Answering,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, 6077–86.
- [4] A. Andreopoulos and J. K. Tsotsos, “50 Years of Object Recognition: Directions Forward,” *Computer Vision and Image Understanding*, 117(8), 827–91.
- [5] A. Bansal, K. Sikka, G. Sharma, R. Chellappa, and A. Divakaran, “Zero-Shot Object Detection,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, 384–400.
- [6] A. Bellet, A. Habrard, and M. Sebban, “A Survey on Metric Learning for Feature Vectors and Structured Data,” *arXiv preprint arXiv:1306.6709*.
- [7] Y. Bengio, A. Courville, and P. Vincent, “Representation Learning: A Review and New Perspectives,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8), 1798–828.
- [8] J. Berant, A. Chou, R. Frostig, and P. Liang, “Semantic Parsing on Freebase from Question-Answer Pairs,” in *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, 2013, 1533–44.

- [9] H. Bilen and A. Vedaldi, “Weakly Supervised Deep Detection Networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, 2846–54.
- [10] A. Borji, M.-M. Cheng, Q. Hou, H. Jiang, and J. Li, “Salient Object Detection: A Survey,” *Computational Visual Media*, 5(2), 117–50.
- [11] Z. Cai and N. Vasconcelos, “Cascade R-CNN: Delving into High Quality Object Detection,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, 6154–62.
- [12] F. Cakir, K. He, X. Xia, B. Kulis, and S. Sclaroff, “Deep metric Learning to Rank,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, 1861–70.
- [13] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-End Object Detection with Transformers,” in *European Conference on Computer Vision*, Springer, 2020, 213–29.
- [14] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, “Emerging Properties in Self-Supervised Vision Transformers,” *arXiv preprint arXiv:2104.14294*.
- [15] Ò. Celma and X. Serra, “FOAFing the Music: Bridging the Semantic Gap in Music Recommendation,” *Journal of Web Semantics*, 6(4), 250–6.
- [16] B. X. Chen and J. K. Tsotsos, “Fast Visual Object Tracking with Rotated Bounding Boxes,” *arXiv preprint arXiv:1907.03892*.
- [17] H.-S. Chen, M. Rouhsedaghat, H. Ghani, S. Hu, S. You, and C.-C. J. Kuo, “DefakeHop: A Light-Weight High-Performance Deepfake Detector,” in *2021 IEEE International Conference on Multimedia and Expo (ICME)*, IEEE, 2021, 1–6.
- [18] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A Simple Framework for Contrastive Learning of Visual Representations,” in *International Conference on Machine Learning*, PMLR, 2020, 1597–607.
- [19] Y. Chen and C.-C. J. Kuo, “PixelHop: A Successive Subspace Learning (SSL) Method for Object Recognition,” *Journal of Visual Communication and Image Representation*, 102749.
- [20] Y. Chen, M. Rouhsedaghat, S. You, R. Rao, and C.-C. J. Kuo, “PixelHop++: A Small Successive-Subspace-Learning-Based (SSL-based) Model for Image Classification,” *arXiv preprint arXiv:2002.03141*.
- [21] Y. Chen, Z. Xu, S. Cai, Y. Lang, and C.-C. J. Kuo, “A Saak Transform Approach to Efficient, Scalable and Robust Handwritten Digits Recognition,” in *2018 Picture Coding Symposium (PCS)*, IEEE, 2018, 174–8.
- [22] Y. Chen, H. Sampathkumar, B. Luo, and X.-w. Chen, “iLike: Bridging the Semantic Gap in Vertical Image Search by Integrating Text and Visual Features,” *IEEE Transactions on Knowledge and Data Engineering*, 25(10), 2257–70.

- [23] S. Chopra, R. Hadsell, and Y. LeCun, "Learning a Similarity Metric Discriminatively, with Application to Face Verification," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, Vol. 1, IEEE, 2005, 539–46.
- [24] R. G. Cinbis, J. Verbeek, and C. Schmid, "Weakly Supervised Object Localization with Multi-Fold Multiple Instance Learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(1), 189–203.
- [25] C. Cortes and V. Vapnik, "Support-Vector Networks," *Machine Learning*, 20(3), 273–97.
- [26] Z. Dai, B. Cai, Y. Lin, and J. Chen, "Up-detr: Unsupervised Pre-training for Object Detection with Transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, 1601–10.
- [27] N. Dalal and B. Triggs, "Histograms of Oriented Gradients for Human Detection," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, Vol. 1, IEEE, 2005, 886–93.
- [28] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-Theoretic Metric Learning," in *Proceedings of the 24th International Conference on Machine Learning*, 2007, 209–16.
- [29] 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*, IEEE, 2009, 248–55.
- [30] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive Angular Margin Loss for Deep Face Recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, 4690–9.
- [31] K. Desai and J. Johnson, "Virtex: Learning Visual Representations from Textual Annotations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, 11162–73.
- [32] P. Dollar, C. Wojek, B. Schiele, and P. Perona, "Pedestrian Detection: An Evaluation of the State of the Art," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(4), 743–61.
- [33] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, *et al.*, "An Image is Worth 16x16 words: Transformers for Image Recognition at Scale," *arXiv preprint arXiv:2010.11929*.
- [34] J. Duan, S. Liao, X. Guo, and S. Z. Li, "Face Detection by Aggregating Visible Components," in *Asian Conference on Computer Vision*, Springer, 2016, 319–33.
- [35] J. Duan, S. Liao, S. Zhou, and S. Z. Li, "Face Classification: A Specialized Benchmark Study," in *Chinese Conference on Biometric Recognition*, Springer, 2016, 22–9.

- [36] J. Duan, Y.-L. Lin, S. Tran, L. S. Davis, and C.-C. J. Kuo, "SLADE: A Self-Training Framework For Distance Metric Learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, 9644–53.
- [37] J. Duan, J. Wan, S. Zhou, X. Guo, and S. Z. Li, "A Unified Framework for Multi-modal Isolated Gesture Recognition," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 14(1s), 1–16.
- [38] J. Duan, S. Zhou, J. Wan, X. Guo, and S. Z. Li, "Multi-modality Fusion Based on Consensus-Voting and 3D Convolution for Isolated Gesture Recognition," *arXiv preprint arXiv:1611.06689*.
- [39] Y. Duan, J. Lu, Z. Wang, J. Feng, and J. Zhou, "Learning Deep Binary Descriptor with Multi-quantization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, 1183–92.
- [40] S. R. Dubey, "A decade Survey of Content based Image Retrieval using Deep Learning," *IEEE Transactions on Circuits and Systems for Video Technology*.
- [41] M. Enzweiler and D. M. Gavrilu, "Monocular Pedestrian Detection: Survey and Experiments," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(12), 2179–95.
- [42] V. Erin Liong, J. Lu, G. Wang, P. Moulin, and J. Zhou, "Deep Hashing for Compact Binary Codes Learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, 2475–83.
- [43] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The Pascal Visual Object Classes (VOC) Challenge," *International Journal of Computer Vision*, 88(2), 303–38.
- [44] P. Felzenszwalb, D. McAllester, and D. Ramanan, "A Discriminatively Trained, Multiscale, Deformable Part Model," in *2008 IEEE Conference on Computer Vision and Pattern Recognition*, IEEE, 2008, 1–8.
- [45] Y. Freund and R. E. Schapire, "A Decision-Theoretic Generalization of On-line Learning and an Application to Boosting," *Journal of Computer and System Sciences*, 55(1), 119–39.
- [46] A. Frome, G. Corrado, J. Shlens, S. Bengio, J. Dean, M. Ranzato, and T. Mikolov, "Devise: A Deep Visual-Semantic Embedding Model."
- [47] J. Gao, J. Wang, S. Dai, L.-J. Li, and R. Nevatia, "Note-RCNN: Noise Tolerant Ensemble RCNN for Semi-supervised Object Detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, 9508–17.
- [48] V. Gattupalli, Y. Zhuo, and B. Li, "Weakly Supervised Deep Image Hashing through Tag Embeddings," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, 10375–84.

- [49] D. Geronimo, A. M. Lopez, A. D. Sappa, and T. Graf, "Survey of Pedestrian Detection for Advanced Driver Assistance Systems," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(7), 1239–58.
- [50] G. Ghiasi, T.-Y. Lin, and Q. V. Le, "Nas-fpn: Learning Scalable Feature Pyramid Architecture for Object Detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, 7036–45.
- [51] R. Girshick, "Fast R-CNN," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, 1440–8.
- [52] 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), 142–58.
- [53] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, 580–7.
- [54] K. Grauman and B. Leibe, "Visual Object Recognition (Synthesis Lectures on Artificial Intelligence and Machine Learning)," *Morgan & Claypool*.
- [55] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai, *et al.*, "Recent Advances in Convolutional Neural Networks," *Pattern Recognition*, 77, 354–77.
- [56] Z. Guan, F. Xie, W. Zhao, X. Wang, L. Chen, W. Zhao, and J. Peng, "Tag-based Weakly-supervised Hashing for Image Retrieval," in *IJCAI*, 2018, 3776–82.
- [57] D. Gupta, A. Anantharaman, N. Mangain, V. N. Balasubramanian, C. Jawahar, *et al.*, "A Multi-space Approach to Zero-shot Object Detection," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, 1209–17.
- [58] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality Reduction by Learning an Invariant Mapping," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, Vol. 2, IEEE, 2006, 1735–42.
- [59] R. M. Haralick and L. G. Shapiro, "Image Segmentation Techniques," *Computer Vision, Graphics, and Image Processing*, 29(1), 100–32.
- [60] J. S. Hare, P. H. Lewis, P. G. Enser, and C. J. Sandom, "Mind the Gap: Another Look at the Problem of the Semantic Gap in Image Retrieval," in *Multimedia Content Analysis, Management, and Retrieval 2006*, Vol. 6073, International Society for Optics and Photonics, 2006, 607309.

- [61] J. S. Hare, P. A. Sinclair, P. H. Lewis, K. Martinez, P. G. Enser, and C. J. Sandom, “Bridging the Semantic Gap in Multimedia Information Retrieval: Top-down and Bottom-up Approaches.”
- [62] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, “Momentum Contrast for Unsupervised Visual Representation Learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, 9729–38.
- [63] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask R-CNN,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, 2961–9.
- [64] K. He, X. Zhang, S. Ren, and J. Sun, “Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9), 1904–16.
- [65] J. Hoffman, S. Guadarrama, E. S. Tzeng, R. Hu, J. Donahue, R. Girshick, T. Darrell, and K. Saenko, “LSDA: Large Scale Detection Through Adaptation,” *Advances in Neural Information Processing Systems*, 27, 3536–44.
- [66] S. C. Hoi, W. Liu, and S.-F. Chang, “Semi-supervised Distance Metric Learning for Collaborative Image Retrieval and Clustering,” *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 6(3), 1–26.
- [67] A. G. Howard, “Some Improvements on Deep Convolutional Neural Network based Image Classification,” *arXiv preprint arXiv:1312.5402*.
- [68] C.-K. Hsieh, L. Yang, Y. Cui, T.-Y. Lin, S. Belongie, and D. Estrin, “Collaborative Metric Learning,” in *Proceedings of the 26th International Conference on World Wide Web*, 2017, 193–201.
- [69] Q. Hu, J. Wu, J. Cheng, L. Wu, and H. Lu, “Pseudo Label Based Unsupervised Deep Discriminative Hashing for Image Retrieval,” in *Proceedings of the 25th ACM International Conference on Multimedia*, 2017, 1584–90.
- [70] C. Huang, C. C. Loy, and X. Tang, “Unsupervised Learning of Discriminative Attributes and Visual Representations,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, 5175–84.
- [71] N. Idrissi, J. Martinez, and D. Aboutajdine, “Bridging the Semantic Gap for Texture-based Image Retrieval and Navigation.,” *Journal of Multimedia*, 4(5).
- [72] R. Jin, S. Wang, and Y. Zhou, “Regularized Distance Metric Learning: Theory and Algorithm.,” in *NIPS*, Vol. 22, Citeseer, 2009, 862–70.
- [73] S. Jin, S. Zhou, Y. Liu, C. Chen, X. Sun, H. Yao, and X.-S. Hua, “Ssah: Semi-supervised Adversarial Deep Hashing with Self-paced Hard Sample Generation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34, No. 07, 2020, 11157–64.

- [74] P. Kadam, M. Zhang, S. Liu, and C.-C. J. Kuo, “R-PointHop: A Green, Accurate and Unsupervised Point Cloud Registration Method,” *arXiv preprint arXiv:2103.08129*.
- [75] P. Kadam, M. Zhang, S. Liu, and C.-C. J. Kuo, “Unsupervised Point Cloud Registration via Salient Points Analysis (SPA),” in *2020 IEEE International Conference on Visual Communications and Image Processing (VCIP)*, IEEE, 2020, 5–8.
- [76] M. Kaya and H. Ş. Bilge, “Deep Metric Learning: A Survey,” *Symmetry*, 11(9), 1066.
- [77] Y. Ke and R. Sukthankar, “PCA-SIFT: A More Distinctive Representation for Local Image Descriptors,” in *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004*. Vol. 2, IEEE, 2004, II–II.
- [78] I. Krasin, T. Duerig, N. Aldrin, V. Ferrari, S. Abu-El-Haija, A. Kuznetsova, H. Rom, J. Uijlings, S. Popov, A. Veit, *et al.*, “Openimages: A Public Dataset for Large-scale Multi-label and Multi-class Image Classification,” *Dataset available from <https://github.com/openimages>*, 2(3), 18.
- [79] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, “3d Object Representations for Fine-grained Categorization,” in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2013, 554–61.
- [80] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet Classification with Deep Convolutional Neural Networks,” *Advances in Neural Information Processing Systems*, 25, 1097–105.
- [81] B. Kulis *et al.*, “Metric Learning: A Survey,” *Foundations and Trends in Machine Learning*, 5(4), 287–364.
- [82] C.-C. J. Kuo, “The CNN as a Guided Multilayer RECOs Transform [Lecture Notes],” *IEEE Signal Processing Magazine*, 34(3), 81–9.
- [83] C.-C. J. Kuo, “Understanding Convolutional Neural Networks with a Mathematical Model,” *Journal of Visual Communication and Image Representation*, 41, 406–13.
- [84] C.-C. J. Kuo and Y. Chen, “On Data-driven Saak Transform,” *Journal of Visual Communication and Image Representation*, 50, 237–46.
- [85] C.-C. J. Kuo, M. Zhang, S. Li, J. Duan, and Y. Chen, “Interpretable Convolutional Neural Networks via Feedforward Design,” *Journal of Visual Communication and Image Representation*.
- [86] Y. LeCun, Y. Bengio, and G. Hinton, “Deep Learning,” *Nature*, 521(7553), 436–44.
- [87] J.-E. Lee, R. Jin, and A. K. Jain, “Rank-based Distance Metric Learning: An Application to Image Retrieval,” in *2008 IEEE Conference on Computer Vision and Pattern Recognition*, IEEE, 2008, 1–8.

- [88] X. Lei, G. Zhao, and C.-C. J. Kuo, "NITES: A Non-parametric Interpretable Texture Synthesis Method," in *2020 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, IEEE, 2020, 1698–706.
- [89] X. Lei, G. Zhao, K. Zhang, and C.-C. J. Kuo, "TGHop: An explainable, Efficient and Lightweight Method for Texture Generation," *arXiv preprint arXiv:2107.04020*.
- [90] B. Li, J. H. Errico, H. Pan, and I. Sezan, "Bridging the Semantic Gap in Sports Video Retrieval and Summarization," *Journal of Visual Communication and Image Representation*, 15(3), 393–424.
- [91] C. Li, C. Deng, N. Li, W. Liu, X. Gao, and D. Tao, "Self-supervised Adversarial Hashing Networks for Cross-modal Retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, 4242–51.
- [92] D. Li and Y. Tian, "Survey and Experimental Study on Metric Learning Methods," *Neural Networks*, 105, 447–62.
- [93] S. Li, H. Zhang, J. Zhang, Y. Ren, and C.-C. J. Kuo, "Box Refinement: Object Proposal Enhancement and Pruning," in *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE, 2017, 979–88.
- [94] Y. Li and X. Zhang, "SiamVGG: Visual Tracking Using Deeper Siamese Networks," *arXiv preprint arXiv:1902.02804*.
- [95] Z. Li, J. Tang, L. Zhang, and J. Yang, "Weakly-supervised Semantic Guided Hashing for Social Image Retrieval," *International Journal of Computer Vision*, 128.
- [96] Z. Li, L. Yao, X. Zhang, X. Wang, S. Kanhere, and H. Zhang, "Zero-shot Object Detection with Textual Descriptions," in *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33, No. 01, 2019, 8690–7.
- [97] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, "Person Re-identification by Local Maximal Occurrence Representation and Metric Learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, 2197–206.
- [98] K. Lin, J. Lu, C.-S. Chen, and J. Zhou, "Learning Compact Binary Descriptors with Unsupervised Deep Neural Networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, 1183–92.
- [99] K. Lin, J. Lu, C.-S. Chen, J. Zhou, and M.-T. Sun, "Unsupervised Deep Learning of Compact Binary Descriptors," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(6), 1501–14.
- [100] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common Objects in Context," in *European Conference on Computer Vision*, Springer, 2014, 740–55.

- [101] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. Van Der Laak, B. Van Ginneken, and C. I. Sánchez, “A Survey on Deep Learning in Medical Image Analysis,” *Medical Image Analysis*, 42, 60–88.
- [102] L. Liu, W. Ouyang, X. Wang, P. Fieguth, J. Chen, X. Liu, and M. Pietikäinen, “Deep Learning for Generic Object Detection: A Survey,” *International Journal of Computer Vision*, 128(2), 261–318.
- [103] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “Ssd: Single Shot Multibox detector,” in *European Conference on Computer Vision*, Springer, 2016, 21–37.
- [104] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, “Sphereface: Deep Hypersphere Embedding for Face Recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, 212–20.
- [105] X. Liu, F. Xing, C. Yang, C.-C. J. Kuo, S. Babu, G. E. Fakhri, T. Jenkins, and J. Woo, “VoxelHop: Successive Subspace Learning for ALS Disease Classification Using Structural MRI,” *arXiv preprint arXiv:2101.05131*.
- [106] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin Transformer: Hierarchical Vision Transformer using Shifted Windows,” *arXiv preprint arXiv:2103.14030*.
- [107] D. G. Lowe, “Object Recognition from Local Scale-invariant Features,” in *Proceedings of the Seventh IEEE International Conference on Computer Vision*, Vol. 2, Ieee, 1999, 1150–7.
- [108] H. Ma, J. Zhu, M. R.-T. Lyu, and I. King, “Bridging the Semantic Gap between Image Contents and Tags,” *IEEE Transactions on Multimedia*, 12(5), 462–73.
- [109] A. Manimaran, T. Ramanathan, S. You, and C.-C. J. Kuo, “Visualization, Discriminability and Applications of Interpretable Saak Features,” *Journal of Visual Communication and Image Representation*, 66, 102699.
- [110] B. McFee and G. R. Lanckriet, “Metric Learning to Rank,” in *ICML*, 2010.
- [111] K. Mikolajczyk and C. Schmid, “A Performance Evaluation of Local Descriptors,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(10), 1615–30.
- [112] I. Misra and L. v. d. Maaten, “Self-supervised Learning of Pretext-invariant Representations,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, 6707–17.
- [113] S. Moran and V. Lavrenko, “Sparse Kernel Learning for Image Annotation,” in *Proceedings of International Conference on Multimedia Retrieval*, 2014, 113–20.

- [114] Y. Movshovitz-Attias, A. Toshev, T. K. Leung, S. Ioffe, and S. Singh, “No Fuss Distance Metric Learning using Proxies,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, 360–8.
- [115] K. Musgrave, S. Belongie, and S.-N. Lim, “A Metric Learning Reality Check,” in *European Conference on Computer Vision*, Springer, 2020, 681–99.
- [116] H. V. Nguyen and L. Bai, “Cosine Similarity Metric Learning for Face Verification,” in *Asian Conference on Computer Vision*, Springer, 2010, 709–20.
- [117] H. Oh Song, Y. Xiang, S. Jegelka, and S. Savarese, “Deep Metric Learning via Lifted Structured Feature Embedding,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, 4004–12.
- [118] T. Ojala, M. Pietikainen, and T. Maenpaa, “Multiresolution Gray-scale and Rotation Invariant Texture Classification with Local Binary Patterns,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(7), 971–87.
- [119] Y. Pang, Y. Li, J. Shen, and L. Shao, “Towards bridging semantic gap to improve semantic segmentation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, 4230–9.
- [120] Papers With Code, *Coco Test-Dev Benchmark (Object Detection)*, <https://paperswithcode.com/sota/object-detection-on-coco>, (accessed: 08.19.2021).
- [121] Papers With Code, *Image Classification on ImageNet*, <https://paperswithcode.com/sota/image-classification-on-imagenet>, (accessed: 08.19.2021).
- [122] M. Paulin, M. Douze, Z. Harchaoui, J. Mairal, F. Perronin, and C. Schmid, “Local Convolutional Features with Unsupervised Training for Image Retrieval,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, 91–9.
- [123] H. Pham, Z. Dai, Q. Xie, and Q. V. Le, “Meta Pseudo Labels,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, 11557–68.
- [124] Q. Qian, L. Shang, B. Sun, J. Hu, H. Li, and R. Jin, “Softtriple Loss: Deep Metric Learning without Triplet Sampling,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, 6450–8.
- [125] S. Rahman, S. Khan, and N. Barnes, “Improved Visual-semantic Alignment for Zero-shot Object Detection,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34, No. 07, 2020, 11932–9.
- [126] S. Rahman, S. Khan, and N. Barnes, “Transductive Learning for Zero-shot Object Detection,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, 6082–91.

- [127] V. Ramanathan, R. Wang, and D. Mahajan, "Dlwl: Improving Detection for Lowshot Classes with Weakly Labelled Data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, 9342–52.
- [128] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-time Object Detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, 779–88.
- [129] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-time Object Detection with Region Proposal Networks," *Advances in Neural Information Processing Systems*, 28, 91–9.
- [130] Y. Ren, C. Zhu, and S. Xiao, "Deformable Faster R-CNN with Aggregating Multi-layer Features for Partially Occluded Object Detection in Optical Remote Sensing Images," *Remote Sensing*, 10(9), 1470.
- [131] M. Rouhsedaghat, M. Monajatipoor, Z. Azizi, and C.-C. J. Kuo, "Successive Subspace Learning: An Overview," *arXiv preprint arXiv:2103.00121*.
- [132] M. Rouhsedaghat, Y. Wang, X. Ge, S. Hu, S. You, and C.-C. J. Kuo, "Facehop: A Light-weight Low-resolution Face Gender Classification Method," *arXiv preprint arXiv:2007.09510*.
- [133] M. Rouhsedaghat, Y. Wang, S. Hu, S. You, and C.-C. J. Kuo, "Low-resolution Face Recognition in Resource-constrained Environments," *Pattern Recognition Letters*.
- [134] Y. Rui, T. S. Huang, and S.-F. Chang, "Image Retrieval: Current Techniques, Promising Directions, and Open Issues," *Journal of Visual Communication and Image Representation*, 10(1), 39–62.
- [135] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, *et al.*, "Imagenet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision*, 115(3), 211–52.
- [136] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "Overfeat: Integrated Recognition, Localization and Detection using Convolutional Networks," *arXiv preprint arXiv:1312.6229*.
- [137] I. K. Sethi, I. L. Coman, and D. Stan, "Mining Association Rules between Low-level Image Features and High-level Concepts," in *Data Mining and Knowledge Discovery: Theory, Tools, and Technology III*, Vol. 4384, International Society for Optics and Photonics, 2001, 279–90.
- [138] M. J. Shafiee, B. Chywl, F. Li, and A. Wong, "Fast YOLO: A Fast You Only Look Once System for Real-time Embedded Object Detection in Video," *arXiv preprint arXiv:1709.05943*.
- [139] A. W. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, "Content-based Image Retrieval at the End of the Early Years," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(12), 1349–80.

- [140] K. Sohn, “Improved Deep Metric Learning with Multi-class N-pair Loss Objective,” in *Advances in Neural Information Processing Systems*, 2016, 1857–65.
- [141] C. Sun, C. Gan, and R. Nevatia, “Automatic Concept Discovery from Parallel Text and Visual Corpora,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, 2596–604.
- [142] Y. Sun, C. Cheng, Y. Zhang, C. Zhang, L. Zheng, Z. Wang, and Y. Wei, “Circle Loss: A Unified Perspective of Pair Similarity Optimization,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, 6398–407.
- [143] Z. Sun, G. Bebis, and R. Miller, “On-road Vehicle Detection: A Review,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(5), 694–711.
- [144] V. Svetnik, A. Liaw, C. Tong, J. C. Culberson, R. P. Sheridan, and B. P. Feuston, “Random Forest: a Classification and Regression Tool for Compound Classification and QSAR Modeling,” *Journal of Chemical Information and Computer Sciences*, 43(6), 1947–58.
- [145] M. J. Swain and D. H. Ballard, “Color indexing,” *International journal of computer vision*, 7(1), 11–32.
- [146] C. Szegedy, A. Toshev, and D. Erhan, “Deep Neural Networks for Object Detection.”
- [147] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the Inception Architecture for Computer Vision,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, 2818–26.
- [148] J. Tang and Z. Li, “Weakly Supervised Multimodal Hashing for Scalable Social Image Retrieval,” *IEEE Transactions on Circuits and Systems for Video Technology*, 28(10), 2730–41.
- [149] Y. Tang, J. Wang, B. Gao, E. Dellandréa, R. Gaizauskas, and L. Chen, “Large Scale Semi-supervised Object Detection using Visual and Semantic Knowledge Transfer,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, 2119–28.
- [150] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, “Training Data-efficient Image Transformers & Distillation through Attention,” in *International Conference on Machine Learning*, PMLR, 2021, 10347–57.
- [151] H. Touvron, A. Vedaldi, M. Douze, and H. Jégou, “Fixing the Train-test Resolution Discrepancy,” *arXiv preprint arXiv:1906.06423*.
- [152] T.-W. Tseng, K.-J. Yang, C.-C. J. Kuo, and S.-H. Tsai, “An Interpretable Compression and Classification System: Theory and Applications,” *IEEE Access*, 8, 143962–74.

- [153] O. Tuzel, F. Porikli, and P. Meer, “Region Covariance: A Fast Descriptor for Detection and Classification,” in *European Conference on Computer Vision*, Springer, 2006, 589–600.
- [154] J. Uijlings, S. Popov, and V. Ferrari, “Revisiting Knowledge Transfer for Training Object Class Detectors,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, 1101–10.
- [155] J. R. Uijlings, K. E. Van De Sande, T. Gevers, and A. W. Smeulders, “Selective Search for Object Recognition,” *International Journal of Computer Vision*, 104(2), 154–71.
- [156] A. van den Oord, Y. Li, and O. Vinyals, “Representation Learning with Contrastive Predictive Coding,” *arXiv preprint arXiv:1807.03748*.
- [157] L. Van der Maaten and G. Hinton, “Visualizing Data Using t-SNE.,” *Journal of Machine Learning Research*, 9(11).
- [158] S. Vembu, M. Kiesel, M. Sintek, and S. Baumann, “Towards Bridging the Semantic Gap in Multimedia Annotation and Retrieval,” in *1st International Workshop on Semantic Web Annotations for Multimedia (SWAMM)*, 2006.
- [159] P. Viola and M. Jones, “Rapid Object Detection using a Boosted Cascade of Simple Features,” in *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, Vol. 1, Ieee, 2001, I–I.
- [160] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, “The Caltech-UCSD Birds-200-2011 Dataset.”
- [161] F. Wan, C. Liu, W. Ke, X. Ji, J. Jiao, and Q. Ye, “C-mil: Continuation Multiple Instance Learning for Weakly Supervised Object Detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, 2199–208.
- [162] F. Wang and J. Sun, “Survey on Distance Metric Learning and Dimensionality Reduction in Data Mining,” *Data Mining and Knowledge Discovery*, 29(2), 534–64.
- [163] F. Wang, J. Cheng, W. Liu, and H. Liu, “Additive Margin Softmax for Face Verification,” *IEEE Signal Processing Letters*, 25(7), 926–30.
- [164] F. Wang, X. Xiang, J. Cheng, and A. L. Yuille, “Normface: L2 Hypersphere Embedding for Face Verification,” in *Proceedings of the 25th ACM International Conference on Multimedia*, 2017, 1041–9.
- [165] G. Wang, Q. Hu, J. Cheng, and Z. Hou, “Semi-supervised Generative Adversarial Hashing for Image Retrieval,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, 469–85.
- [166] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, “Cosface: Large Margin Cosine Loss for Deep Face Recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, 5265–74.

- [167] J. Wang, F. Zhou, S. Wen, X. Liu, and Y. Lin, “Deep Metric Learning with Angular Loss,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, 2593–601.
- [168] Y.-X. Wang and M. Hebert, “Model Recommendation: Generating Object Detectors from Few Samples,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, 1619–28.
- [169] X. Wang, X. Han, W. Huang, D. Dong, and M. R. Scott, “Multi-similarity Loss with General Pair Weighting for Deep Metric Learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, 5022–30.
- [170] K. Q. Weinberger and L. K. Saul, “Distance Metric Learning for Large Margin Nearest Neighbor Classification.,” *Journal of Machine Learning Research*, 10(2).
- [171] C.-Y. Wu, R. Manmatha, A. J. Smola, and P. Krahenbuhl, “Sampling Matters in Deep Embedding Learning,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, 2840–8.
- [172] F. Xiong, M. Gou, O. Camps, and M. Szaier, “Person Re-identification using Kernel-based Metric Learning Methods,” in *European Conference on Computer Vision*, Springer, 2014, 1–16.
- [173] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei, “Scene Graph Generation by Iterative Message Passing,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, 5410–9.
- [174] Z. Xu, B. Li, M. Geng, Y. Yuan, and G. Yu, “AnchorFace: An Anchor-based Facial Landmark Detector Across Large Poses,” *arXiv preprint arXiv:2007.03221*.
- [175] B. Yan, H. Peng, J. Fu, D. Wang, and H. Lu, “Learning Spatio-temporal Transformer for Visual Tracking,” *arXiv preprint arXiv:2103.17154*.
- [176] L. Yang, “An Overview of Distance Metric Learning,” in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2007.
- [177] L. Yang and R. Jin, “Distance Metric Learning: A Comprehensive Survey,” *Michigan State University*, 2(2), 4.
- [178] L. Yang, R. Jin, L. Mummert, R. Sukthankar, A. Goode, B. Zheng, S. C. Hoi, and M. Satyanarayanan, “A Boosting Framework for Visuality-preserving Distance Metric Learning and Its Application to Medical Image Retrieval,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(1), 30–44.
- [179] M.-H. Yang, D. J. Kriegman, and N. Ahuja, “Detecting Faces in Images: A Survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(1), 34–58.
- [180] Y. Yang, V. Magouliantitis, and C.-C. J. Kuo, “E-PixelHop: An Enhanced PixelHop Method for Object Classification,” *arXiv preprint arXiv:2107.02966*.

- [181] B. Yao and L. Fei-Fei, “Modeling Mutual Context of Object and Human Pose in Human-object Interaction Activities,” in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, IEEE, 2010, 17–24.
- [182] K. Ye, M. Zhang, A. Kovashka, W. Li, D. Qin, and J. Berent, “Cap2det: Learning to Amplify Weak Caption Supervision for Object Detection,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, 9686–95.
- [183] K. Ye, M. Zhang, W. Li, D. Qin, A. Kovashka, and J. Berent, “Learning to Discover and Localize Visual Objects with Open Vocabulary,” *arXiv preprint arXiv:1811.10080*.
- [184] Q. Ye and D. Doermann, “Text Detection and Recognition in Imagery: A Survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(7), 1480–500.
- [185] D. Yi, Z. Lei, S. Liao, and S. Z. Li, “Deep Metric Learning for Person Re-identification,” in *2014 22nd International Conference on Pattern Recognition*, IEEE, 2014, 34–9.
- [186] B. Yu and D. Tao, “Deep Metric Learning with Tuplet Margin Loss,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, 6490–9.
- [187] A. Zadeh, Y. Chong Lim, T. Baltrusaitis, and L.-P. Morency, “Convolutional Experts Constrained Local Model for 3D Facial Landmark Detection,” in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, 2519–28.
- [188] S. Zafeiriou, C. Zhang, and Z. Zhang, “A Survey on Face Detection in the Wild: Past, Present and Future,” *Computer Vision and Image Understanding*, 138, 1–24.
- [189] A. Zareian, K. D. Rosa, D. H. Hu, and S.-F. Chang, “Open-vocabulary Object Detection using Captions,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, 14393–402.
- [190] A. Zhai and H.-Y. Wu, “Classification is a Strong Baseline for Deep Metric Learning,” *arXiv preprint arXiv:1811.12649*.
- [191] H. Zhang, M. Wang, R. Hong, and T.-S. Chua, “Play and Rewind: Optimizing Binary Representations of Videos by Self-supervised Temporal Hashing,” in *Proceedings of the 24th ACM International Conference on Multimedia*, 2016, 781–90.
- [192] J. Zhang and Y. Peng, “SSDH: Semi-supervised Deep Hashing for Large Scale Image Retrieval,” *IEEE Transactions on Circuits and Systems for Video Technology*, 29(1), 212–25.
- [193] K. Zhang, B. Wang, W. Wang, F. Sohrab, M. Gabbouj, and C.-C. J. Kuo, “AnomalyHop: An SSL-based Image Anomaly Localization Method,” *arXiv preprint arXiv:2105.03797*.

- [194] M. Zhang, P. Kadam, S. Liu, and C.-C. J. Kuo, “Unsupervised Feed-forward Feature (UFF) Learning for Point Cloud Classification and Segmentation,” in *2020 IEEE International Conference on Visual Communications and Image Processing (VCIP)*, IEEE, 2020, 144–7.
- [195] M. Zhang, Y. Wang, P. Kadam, S. Liu, and C.-C. J. Kuo, “PointHop++: A Lightweight Learning Model on Point Sets for 3D Classification,” in *2020 IEEE International Conference on Image Processing (ICIP)*, IEEE, 2020, 3319–23.
- [196] M. Zhang, H. You, P. Kadam, S. Liu, and C.-C. J. Kuo, “PointHop: An Explainable Machine Learning Method for Point Cloud Classification,” *IEEE Transactions on Multimedia*.
- [197] S. Zhang, J. Li, and B. Zhang, “Pairwise Teacher-student Network for Semi-supervised Hashing,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [198] X. Zhang, Y.-H. Yang, Z. Han, H. Wang, and C. Gao, “Object Class Detection: A Survey,” *ACM Computing Surveys (CSUR)*, 46(1), 1–53.
- [199] Z. Zhang, H. Peng, J. Fu, B. Li, and W. Hu, “Ocean: Object-aware Anchor-free Tracking,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16*, Springer, 2020, 771–87.
- [200] R. Zhao and W. I. Grosky, “Bridging the Semantic Gap in Image Retrieval,” in *Distributed Multimedia Databases: Techniques and Applications*, IGI Global, 2002, 14–36.
- [201] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, “Scalable Person Re-identification: A Benchmark,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, 1116–24.
- [202] P. Zhu, H. Wang, and V. Saligrama, “Don’t Even Look Once: Synthesizing Features for Zero-Shot Detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, 11693–702.
- [203] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, “Deformable DETR: Deformable Transformers for End-to-end Object Detection,” *arXiv preprint arXiv:2010.04159*.
- [204] C. L. Zitnick and P. Dollár, “Edge Boxes: Locating Object Proposals from Edges,” in *European Conference on Computer Vision*, Springer, 2014, 391–405.
- [205] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, “Learning Transferable Architectures for Scalable Image Recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, 8697–710.