# A Selective Survey on Versatile Knowledge Distillation Paradigm for Neural Network Models

Jeong-Hoe Ku, JiHun Oh, YoungYoon Lee, Gaurav Pooniwala, SangJeong Lee

*Samsung Research,*

*Samsung Electronics Co., Ltd.*

Seoul, Republic of Korea

Email: {mrku69, jihun2331.oh, euler.lee, pt.gaurav, sj94.lee}@samsung.com

*Abstract*—**This paper aims to provide a selective survey about *knowledge distillation* (KD) framework for researchers and practitioners to take advantage of it for developing new optimized models in the deep neural network field. To this end, we give a brief overview of knowledge distillation and some related works including *learning using privileged information* (LUPI) and *generalized distillation* (GD). Even though knowledge distillation based on the teacher-student architecture was initially devised as a model compression technique, it has found versatile applications over various frameworks. In this paper, we review the characteristics of knowledge distillation from the hypothesis that the three important ingredients of knowledge distillation are *distilled knowledge and loss*, *teacher-student paradigm*, and the *distillation process*. In addition, we survey the versatility of the knowledge distillation by studying its direct *applications* and its usage in combination with other *deep learning paradigms*. Finally we present some future works in knowledge distillation including explainable knowledge distillation where the analytical analysis of the performance gain is studied and the self-supervised learning which is a hot research topic in deep learning community.**

*Index Terms*—**Knowledge Distillation, Compression, Quantization, Object Detection**

## I. Introduction

Deep neural networks have been applied to various tasks and achieved dramatic success in many fields. However, larger neural networks with more layers and nodes are used to achieve higher performance. Many recent successful deep neural models are computationally expensive and memory intensive [1]. Since it is difficult to deploy such heavy models on devices with low system resources for running real-time applications with a strict latency requirement, it is crucial to reduce their parameter size and computational complexity without performance degradation.

Yu Cheng *et al.* [2] briefly categorized compression techniques for the purpose into four schemes: *parameter pruning and sharing, low-rank factorization, transferred/compact convolution filters, quantization*, and *knowledge distillation*, each of which has it's own advantages and drawbacks. Among these techniques, we focus on ***knowledge distillation*** as it is an empirically very successful technique for knowledge transfer between classifiers in an interactive manner which is more similar to how humans learn. Knowledge distillation with neural networks was pioneered by Hinton *et al.* [6] which is a transfer learning method that aims to improve the training of a student network by relying on knowledge borrowed from a powerful teacher network.

This paper is organized as follows:
In Section II, we explain knowledge distillation framework briefly and reviewed three important design factors of knowledge distillation including Distilled Knowledge and Loss, the Teacher-Student architecture and the Distillation Process to investigate the internals of knowledge distillation framework. In Section III, we survey versatile knowledge distillation framework usages in the fields of Computer Vision, Natural Language Processing and Quantization. In Section IV, we look at the position of Knowledge distillation in the broader Deep learning paradigms of Supervised Learning, Weakly-Supervised Learning and Semi-Supervised/Unsupervised Learning. In Section V we look at a few techniques that do not fit the definition of Knowledge distillation but use very similar ideas. In Section VI, we suggest two topics as future works of knowledge distillation

Fig. 16, Fig. 17, Fig. 18, and Fig. 19 in Annex show the brief overview of the topics of knowledge distillation reviewed in this survey paper.

## II. Knowledge Distillation

Knowledge distillation is an effective model compression technique in which a compact model (student) is trained under the supervision of a larger pre-trained model or an ensemble of models (teacher).

Knowledge distillation aims to improve the performance of the student network by providing additional supervision from a teacher network. To the best of our knowledge, exploiting knowledge transfer to compress model was first proposed in C. Buciluă *et al.* [4]. They trained a compressed/ensemble model of strong classifiers with pseudo-labeled data, and reproduced the output of the original larger network. However, the work is limited to shallow models. The idea has been adopted in [5] as knowledge distillation to compress deep and wide networks into shallower ones, where the compressed model mimicked the function learned by the complex model. Hinton *et al.* [6] popularized the concept of Knowledge Distillation to be extended to more practical uses. The work in [6] proposed knowledge distillation as a more general case of C. Buciluă *et al.* [4] by adopting the concept of temperature parameter at

the output of teacher. The student was trained to predict the output and the classification labels.

The main idea of knowledge distillation approach is to shift knowledge from a large teacher model into a small one by learning the class distributions output via softmax [6]. It has even been observed that the student learns much faster and more reliably if trained using outputs of teacher as soft labels, instead of one-hot-encoded labels.

Since then, a number of knowledge distillation methods have been proposed, each trying to capture and transfer some characteristics of the teacher such as the representation space, decision boundary or intra-data relationship. Despite its simplicity, knowledge distillation demonstrates promising results in various image classification tasks.

Knowledge distillation has proven empirically to be an effective technique for training a compact model [7], [11], [17].

### A. Distilled Knowledge and Loss

Despite the recent advances of knowledge distillation technique, a clear understanding of where knowledge resides in a deep neural network and an optimal method for capturing knowledge from teacher and transferring it to student remains an open question.

In recent advances of knowledge distillation, many forms of knowledge have been defined (Jiaxi Tang *et al.*, 2020 [18]) based on the teacher-student learning paradigm and have shown dramatic success and were analyzed empirically:

- Layer activation [7]
- Auxiliary information [23]
- Jacobian matrix of the model parameters [72], [73]
- Gram matrix derived from pairs of layers [74]
- Activation boundary [34]

Distillation loss for knowledge distillation training is a key factor which is used to penalize the student to transfer this Knowledge from the Teacher to the Student.

Fahad Sarfraz *et al.* [75] presented broad categorization of a diverse set of knowledge distillation methods which differ from each other with respect to how knowledge is defined and transferred from the teacher. Borrowing from their categorization, we cite two groups below to demonstrate how to capture the knowledge from teacher.

**a) Response Distillation** uses only the outputs of a Teacher to train the student to mimic it. C. Buciluă *et al.* [4] proposed to use the logits of a teacher network as target for the student and to minimize the squared difference. Hinton *et al.* [6] proposed to minimize the KL divergence between the smoother output probabilities. In the original formulation, Hinton *et al.* [6] introduced a knowledge distillation compression framework and proposed mimicking the softened softmax output of the teacher using a temperature parameter. It raised the temperature of the final softmax function and minimize the Kullback-Leibler (KL) divergence between the smoother output probabilities. This softened output transfers more important information which is called *dark knowledge* compared to the hard output. When the soft targets have high entropy, they provide much more information per training case than hard targets and much less variance in the gradient between training cases, so the small model can often be trained on much less data than the original cumbersome model while using a much higher learning rate.

When the dimension of both outputs are the same, these methods can be applied to any pair of network architectures. Although this loss was originally proposed to apply to a simple task such as image classification, it has seen a wide variety of applications.

**b) Representation Space Distillation** aims to mimic the latent feature space of the teacher. Adriana Romero *et al.* [7] introduced intermediate level hints from the teacher's hidden layers to guide the training process of the student. Due to the differences of a dimension of hidden layers between the teacher and student network, design of the feature distillation method needs to be done carefully to prevent information loss when transferring knowledge. Byeong Heo *et al.* [34] proposed a novel feature distillation method and designed a new distillation loss to minimize the information loss. This gave a hint to extend knowledge distillation to more complex tasks such as object detection.

### B. Teacher-Student Architecture

Knowledge distillation has proven to be an effective technique for training a compact model and also providing greater architectural flexibility since it allows for structural differences in the teacher and student networks. There are several variants of the conventional knowledge distillation Student and Teacher architectures which extend the student-teacher learning paradigm to improve the performance and overcome some weak points.

#### 1) Single Teacher- Single Student

In the perspective of a Single Teacher- Single Student learning paradigm, knowledge distillation is a simple way to transfer knowledge of a teacher to improve the performance of small deep learning model called a student. More specifically, knowledge distillation refers to the method that helps the training process of a small network (student) under the supervision of a large network (teacher). The additional supervision about the relative probabilities of secondary class and relational information between data points at the output of Teacher network can be useful in increasing the efficacy of the Student network.

We can downsize a student network regardless of the structural difference between teacher and student. Allowing this architectural flexibility, knowledge distillation is emerging as a next generation approach of network compression. However, too excessive gap of the capacity between a teacher network and a student network is a critical obstacle for knowledge transfer performance. This is empirically shown and analyzed in [17].

#### 2) Multi-Step Learning

Seyed-Iman Mirzadeh *et al.* [17] showed that the student network performance degrades when the gap between student

and teacher networks is large. So, given a fixed student network, one cannot employ an arbitrarily large teacher network; in other words, a teacher network can effectively transfer its knowledge to student networks having up to a certain capacity.

To alleviate this shortcoming, multi-step knowledge distillation was introduced. It is a new distillation framework called *Teacher Assistant Knowledge Distillation* (TAKD), which introduces intermediate-sized network known as teacher assistants (TAs) between the teacher and the student to fill in the gap. TA models are distilled from the teacher, and the student is then only distilled from the TAs. Through extensive empirical evaluations and a theoretical justification, they showed that introducing intermediate TA networks improve the distillation performance and concluded that the size (capacity) gap difference between a teacher-TA and a TA-student is important.
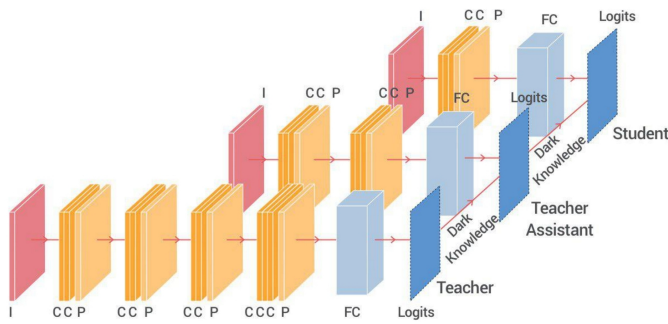


Fig. 1. A teacher assistant network fills the gap between student and teacher networks [17].

Fig. 1 shows the overall knowledge distillation structure incorporating teacher assistant.

*3) Multiple-Teacher Learning*

Shan You *et al.* [15] proposed a new method which uses multiple teacher networks to train a thin and deep student network. Fig. 2 shows overall knowledge distillation incorporating multiple teachers. The method uses three losses to train
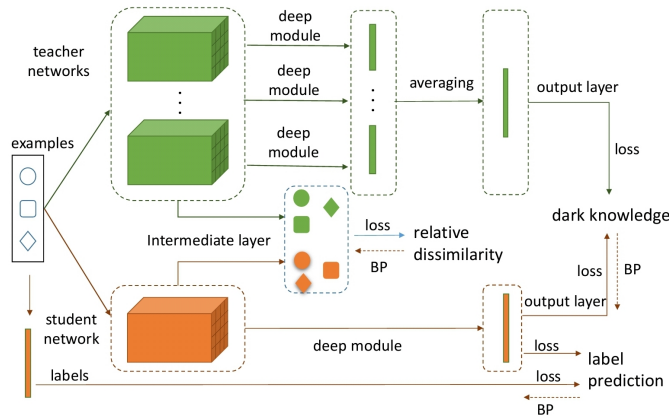


Fig. 2. A graphical diagram for the proposed method to train a new thin deep student network by incorporating multiple comparable teacher networks [15].

the Student: a label prediction loss, a dark knowledge loss and

a relative similarity loss. The incorporation of multiple teacher networks exists in two places.

One is in the output layers via averaging the softened output targets(dark knowledge) from different Teacher networks and using it to train the Student.

It also used knowledge about the intermediate layers by imposing a constraint of the dissimilarity among examples. The authors suggest that relative dissimilarity between intermediate representations of different examples serves as more flexible and appropriate guidance from teacher networks.

### C. Distillation Process

#### 1) Off-line Distillation

In vanilla knowledge distillation (Hinton et al., [6]), we start with a powerful large and pre-trained teacher network and perform one-way knowledge transfer to a small untrained student. This is known as offline knowledge distillation. Most previous Knowledge distillation methods use this process. The large teacher model is first trained on a set of training samples. The teacher model is then used to extract the knowledge in the forms of logits or the intermediate features, which are then used to guide the training of the student model during distillation.

The teacher models typically need to have a high capacity and require a lot of time and data for training. The training of the student model in offline distillation is usually efficient under the guidance of the teacher model. A capacity gap between large teacher network and small student network always exists and Seyed-Iman Mirzadeh *et al.* [17] showed that the student network performance degrades when the gap is too large.
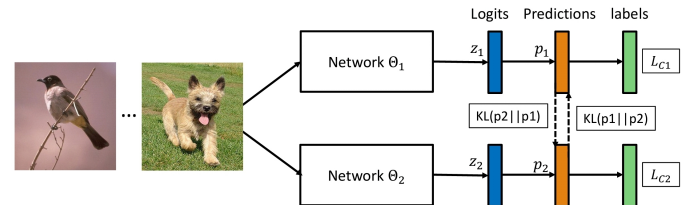
#### 2) On-line Distillation



Fig. 3. The Deep Mutual Learning (DML) schematic [16].

Online distillation is proposed to further improve the performance of the student model, especially when a large-capacity high performance teacher model is not available (Zhang et al. [16], ; Chen et al., 2020c). In online distillation, both the teacher model and the student model are updated simultaneously, and the whole knowledge distillation framework is end-to-end trainable

Ying Zhang et al. [16] first presented an online distillation strategy - deep mutual learning (DML). Rather than one way transfer between a static pre-defined teacher network and a student network, an ensemble of student networks learns collaboratively and teaches each other throughout the training process. DML starts with a pool of untrained students who learn simultaneously to solve the task together. Specifically,

each student is trained using two losses: a conventional supervised learning loss and a mimicry loss that aligns each student's class posterior with the class probabilities of other students. Fig. 3 shows the deep mutual learning method. Each network uses for training a supervised learning loss and a KLD-based mimicry loss to match the probability estimates of its peers.

Guo *et al.* [76] improved the generalization ability of this procedure by using an ensemble of soft logits. Chen *et al.* [77] performed two-level distillation during training with multiple auxiliary peers and one group leader to form a diverse set of peer models.

Kim *et al.* [11] used a feature fusion module to construct the teacher classifier. An ensemble of sub-network classifiers transfers its knowledge to the fused classifier and then the fused classifier delivers its knowledge back to each sub-network, mutually teaching one another in an online-knowledge distillation manner.

## III. KD APPLICATIONS

Knowledge distillation is a flexible approach that can be leveraged in many conventional techniques and applications. In this paper we selected three topics where knowledge distillation is actively applied.

### A. KD in Computer Vision (CV)

An excellent example of the successes in computer vision related deep learning can be illustrated with the ImageNet Challenge [64]. This challenge is a contest involving two different components, image classification and object detection tasks.

#### 1) KD in Image Classification

Since knowledge distillation was first introduced in [4], the initial application was the fundamental computer vision task of image classification [5], [6] and demonstrated excellent improvements. The pioneering study on knowledge distillation [6] showed that a shallow or compressed model trained to mimic the behavior of a deeper or more complex model can recover some or all of the accuracy drop

#### 2) KD in Object Detection

Most modern object detection methods make predictions relative to some initial guesses. Two-stage detectors [46], [47] predict boxes w.r.t. proposals, whereas single-stage methods make predictions w.r.t. anchors [48] or a grid of possible object centers [49], [50].

Applying knowledge distillation techniques to multi-class object detection is challenging for several reasons. First, the performance of detection models suffers from more degradation after compression since detection labels are more expensive and thereby usually less voluminous. Second, knowledge distillation is proposed for classification assuming each class is equally important, whereas that is not the case for detection where the background class is far more prevalent. Third, detection is a more complex task that combines elements of both classification and bounding box regression.

We review four papers [12]–[14], [21] where knowledge distillation was applied to object detection.

**a) Global Feature based Object Detection:** It is difficult to apply vanilla knowledge distillation architecture [6] to object detection because it only transfers soft-target that is the logit of penultimate layer. It is necessary to transfer more knowledge for object detection task to the student network. Romero *et al.* [7] first proposed the framework called FitNet where the student network mimics the full feature maps of the teacher network. The intermediate representations of teacher networks are called 'Hint'. FitNet opened a way to apply knowledge distillation to object detection task and since then researchers proposed many forms of object detection frameworks using knowledge distillation.
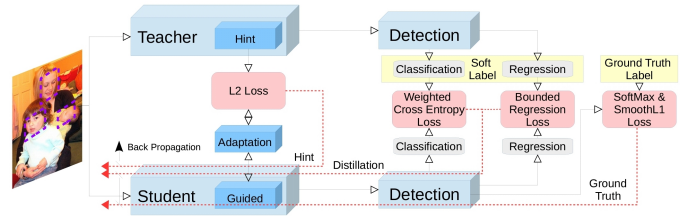


Fig. 4. The proposed learning scheme on a visual object detection task using Faster-RCNN, which mainly consists of the region proposal network (RPN) and the region classification network (RCN) [12].

Guobin Chen *et al.* [12] proposed a new end-to-end trainable framework to train compact and fast multi-class object detection networks with improved accuracy using knowledge distillation [6] and hint learning [7]. Two losses are proposed to effectively address the aforementioned challenges. One is a weighted cross entropy loss for classification that accounts for the imbalance in the impact of misclassification for background class as opposed to object classes. The other is a teacher bounded regression loss for knowledge distillation. For hint learning, adaptation layers are provided to allow the student to better learn from the distribution of neurons in intermediate layers of the teacher. Fig. 4 shows the specialized knowledge distillation scheme proposed for applying to Faster R-CNN. The two networks both use the multi-task loss to jointly learn the classifier and the bounding-box regressor. We employ the final output of the teacher model's RPN and RCN as the distillation targets, and apply the intermediate layer outputs as a hint. Red arrows indicate the backpropagation pathways.

Wanwei Wang *et al.* [14] proposed a clean and effective knowledge distillation method for single-stage object detection called GAN-KD. The feature maps generated by a teacher network and a student network are used as true samples and fake samples respectively, and it conducts adversarial training for both to improve the performance of the student network in single-stage object detection. The GAN algorithm is employed to complete the migration from a teacher network to a student network. GAN-KD is composed of four modules: a teacher network, a student network, a discriminative network (D-Net), and a SSD-Head network as shown below. D-Net learns to determine whether a sample is from the teacher network or the student network. Fig. 5 shows the GAN-KD network
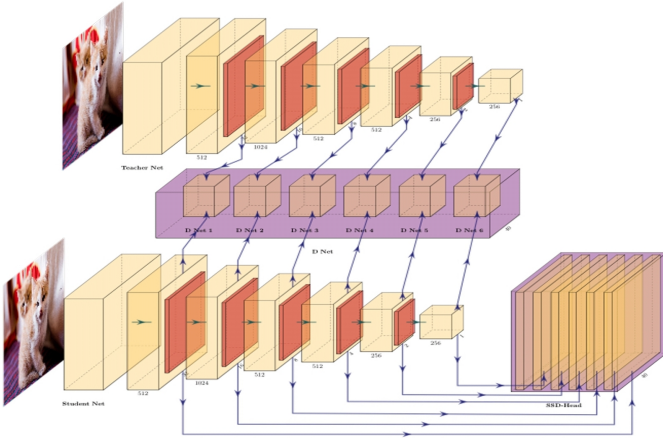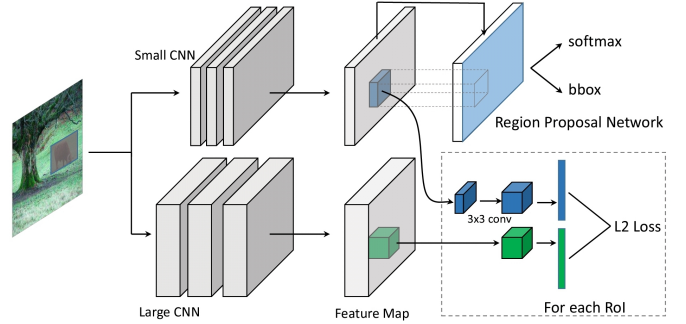
Fig. 5. The GAN-KD network architecture [14].



Fig. 6. An overall architecture of feature mimicking by proposal sampling [13].

on the output structure and internal information of model itself, so it is not applicable to one-stage detector.

architecture. Teacher Net (the top left) is a backbone network of the larger and fully trained SSD model. Student Net (the bottom left) is a smaller network such as MobileNet. SSD-Head (the bottom right) is the head network of the SSD model which is responsible for the classification and anchor box regression. D-Net is a module consisting of six small discriminant networks.

Unlike the traditional knowledge distillation algorithm, there is no need to manually specify the location of knowledge distillation. It does not require the design of complex cost functions, and can be applied to one-stage object detection.

**b) Local Feature based Object Detection:** When applying the whole feature map, the object correspondence information might be ignored or degraded especially for small object feature learning. Compared to the global context features, the features of local regions contain more representative information for object detection. Therefore, Byungseok Roh *et al.* [54] proposed objectness-aware object detection method. There have been some knowledge distillation approaches to leverage local feature based object detection utilizing local features when distilling knowledge from a teacher network to a student network effectively.

Quanquan Li *et al.* [13] presented a feature map mimicking method aiming to train the small model to mimic the feature map activations of the large model in an unified fully convolutional network object detection pipeline. The feature map matters in object detection since both the objectness scores and locations are predicted based on the feature map. Therefore, it is more reasonable to mimic the output feature maps between the two detection networks which contain the response information across an entire image.

Fig. 6 shows the overall architecture based on the RoI-aware knowledge distillation. A Region Proposal Network generates candidate ROIs, which then are used to extract local features from the feature maps.

The work in [13] proposed to only transfer knowledge inside the area of proposals. However, the mimicking regions depend
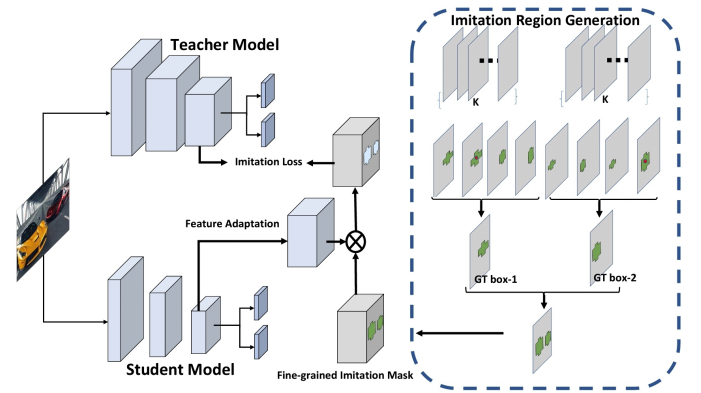


Fig. 7. An illustration of the proposed fine-grained feature imitation method [21].

Tao Wang *et al.* [21] proposed a fine-grained feature imitation method exploiting the cross-location discrepancy of feature response based on the intuition that detectors have attention to the local regions near the object regions. Thus the discrepancy of feature response on the near object anchor locations reveals important information of how a teacher model learns. The novel mechanism is designed to estimate those locations and let student model imitate the teacher on them to achieve the enhanced accuracy. Fig. 7 shows the proposed fine-grained feature imitation method. The student detector is trained by both ground-truth supervision and imitation of teacher's feature response on close object anchor locations. The feature-adaptation layer makes student's guided feature layer compatible with the teacher. To identify informative locations, it iteratively calculates IOU map of each ground-truth bounding box with anchor priors, filter and combine candidates, and generate the final imitation mask.

Recently, Yongcheng Liu *et al.* [55] applied ROI-aware distillation approach to the weakly-supervised detection problem.

*B. Natural Language Processing*

The up-to-date language model, such as BERT, has significantly improved the performance of many natural language

processing. However, the pre-trained language models are computationally expensive and memory intensive, so it is difficult to effectively execute them on resource-restricted devices.
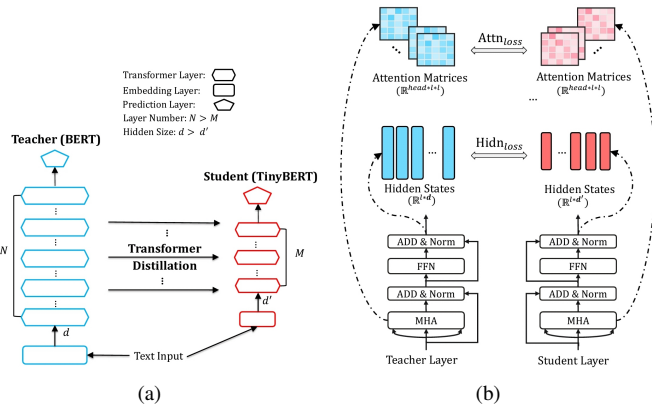
### 1) TinyBERT



Fig. 8. **An overview of Transformer distillation:** (a) the framework of Transformer distillation, (b) the details of Transformer-layer distillation consisting of $Attn_{loss}$ (attention based distillation) and $Hidn_{loss}$ (hidden state based distillation) [56].

Xiaoqi Jiao *et al.* [56] proposed a novel Transformer distillation method that is specially designed for knowledge distillation of the Transformer-based models. The plenty of knowledge encoded in a large "teacher" BERT can be well transferred to a small "student" TinyBERT. They also introduced a new two-stage learning framework for TinyBERT, which performs Transformer distillation at both the pre-training and task-specific learning stages. This framework ensures that TinyBERT can capture the general-domain as well as the task-specific knowledge in BERT. Fig. 8 shows an overview of Transformer distillation.

### 2) DistilBERT

Victor Sanh *et al.* [57] proposed a method to pre-train a smaller general-purpose language representation model, called DistillBERT, which can then be fine-tuned with good performances on a wide range of tasks. Here, they leverage knowledge distillation during the pre-training phase and show that it is possible to reduce the size of a BERT model by 40%, while retaining 97% of its language understanding capabilities and accelerating by 60%. To leverage the inductive biases learned by larger models during pre-training, they introduced a triple loss combining language modeling, distillation, and cosine-distance losses.

### C. Quantization

In the context of deep learning, knowledge distillation has been successfully used to compress heavy networks with a larger capacity model (teacher) to a smaller neural network (student). Quantization is also widely used to reduce parameter size and computational complexity of deep neural networks, especially for resource-constrained edge devices.

In this paper, we review three papers [9]–[11] to survey the cases where knowledge distillation is used for model compression via quantization. In the domain of knowledge distillation, there are a couple of distinctive approaches, based on "*offline knowledge distillation*" and "*online knowledge distillation*". Based on this categorization (Fahad Sarfraz *et al.* [75]), we divided the surveyed cases into the two categories:

- Quantization with KD (Offline KD)
- Quantization-aware KD (Online KD)

### 1) Quantization with KD

In the first approach, a teacher network is fixed when training the quantized student network in the teacher-student network architecture.

Asit Mishra *et al.* [9] studied the combination of quantization and KD, and showed that the performance of low-precision networks can be significantly improved by using knowledge distillation techniques. Its approach, Apprentice, achieves state-of-the-art accuracies using ternary precision and 4-bit precision for variants of ResNet architecture on ImageNet dataset. It presents three schemes on how we can apply knowledge distillation techniques to various stages of the train-and-deploy pipeline. In the first scheme, a low-precision network and a full-precision network are jointly trained from scratch using the knowledge distillation scheme. In the second scheme, it start with a full-precision trained network and transfer knowledge from this trained network continuously to train a low-precision network from scratch. In the third scheme, it starts with a trained full-precision large network and an apprentice network that has been initialized with full-precision weights.
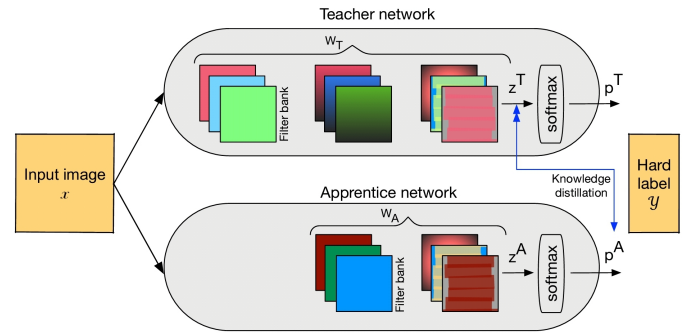


Fig. 9. The schematic structure of the knowledge distillation setup [9].

The apprentice network's precision is lowered and is fine-tuned using knowledge distillation techniques. Each of the scheme produces a low-precision model that surpasses the accuracy of the equivalent low-precision model published to date. Fig. 9 shows the knowledge distillation setup for the apprentice network. Herein, we can regard the high precision network as the teacher network and the low-precision network as the apprentice network.

Antonio Polino *et al.* [10] examined whether distillation and quantization can be jointly leveraged for better compression and proposed two new compression techniques (*quantized distillation* and *differentiable quantization*), which jointly lever-

---

**Algorithm 1:** Quantized Distillation

**Input:** the network weights $w$, quantization level $s$
**Output:** $w^q$

**1 while do**

**2**     $w^q \leftarrow$ quant-function$(w, s)$;

**3**     Run forward pass and compute distillation loss $l(w^q)$ ;

**4**     Run backward pass and compute $\frac{\partial l(w^q)}{\partial w^q}$ ;

**5**     Update original weights using SGD **in full precision** $w = w - \nu \cdot \frac{\partial l(w^q)}{\partial w^q}$

**6 end**

**7** Finally quantize the weights before returning: $w^q \leftarrow$ quant-function$(w, s)$;
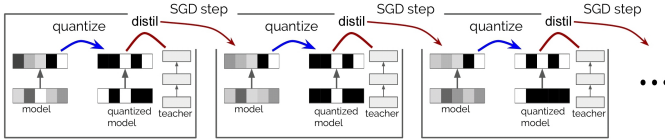
---



Fig. 10. A depiction of the steps of quantized distillation [10].

age weight quantization and distillation on larger networks, called "teacher", into compressed "student" networks. The first method is called *quantized distillation* which leverages distillation during the training process, by incorporating distillation loss, expressed with respect to the teacher network, into the training of a smaller student network whose weights are quantized to a limited set of levels. In other words, it performs the stochastic gradient descent (SGD) step on the full-precision model, but computing the gradient on the quantized model, derived with respect to the distillation loss. Algorithm 1 describes the quantization distillation procedure and Fig. 10 shows quantization distillation steps, respectively. Note the accumulation over multiple steps of gradients on the unquantized model leads to a switch in quantization (e.q. top layer left most square)

---

**Algorithm 2:** Differentiable Quantization

**Input:** the network weights $w$, the initial quantization points $p$
**Output:** $w^q$

**1 while do**

**2**     $w^q \leftarrow$ quant-function$(w, p)$;

**3**     Run forward pass and compute distillation loss $l(w^q)$ ;

**4**     Run backward pass and compute $\frac{\partial l(w^q)}{\partial w^q}$ ;

**5**     Compute [Q(): Quantization function]
$$\frac{\partial Q(v,p)_i}{\partial p_j} = \begin{cases} \alpha_i, & \text{if } v_i \text{ is quantized to } p_j \\ 0, & \text{otherwise} \end{cases} ;$$

**6**     Update original weights using SGD or similar: $p = p - \nu \cdot \frac{\partial l(w^q)}{\partial p}$ ;

**7 end**

---

The second method, *differentiable quantization*, optimizes the location of quantization points through SGD, to better fit the behavior of the teacher model. It is introduced as a general method of improving the accuracy of a quantized neural network, by exploiting non-uniform quantization point placement. Algorithm 2 shows the differentiable quantization procedure.

It shows that quantized shallow students can reach similar accuracy levels to state-of-the-art full precision teacher models, while providing up to the order of magnitude compression, and the inference speed-up factor that is almost linear to the depth reduction.

*2) Quantization-Aware KD*

The second approach adopts online knowledge distillation where a teacher network is also trained when training a quantized student network.

The inherent differences between the distributions of the full-precision teacher network and the low-precision student network may yield difficulty in knowledge transferring from teacher network to student network [16], [17]. Based on the assumption, several studies hypothesized that using a fixed teacher can limit the knowledge transfer [11]. In specific, they argued that using a fixed teacher, as in [9], [10] can limit the knowledge transfer due to the inherent differences between the distributions of the full-precision teacher model and the low-precision student network. Jangho Kim *et al.* [11] tackles this problem via online co-studying (CS) and offline tutoring (TU).
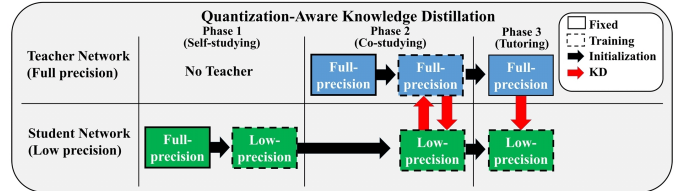


Fig. 11. The overall process of QKD [11].

Jangho Kim *et al.* [11] proposed Quantization-aware Knowledge Distillation (QKD) wherein quantization and KD are carefully coordinated in three phases. First, Self-studying (SS) phase fine-tunes a quantized low-precision student network without KD to obtain a good initialization, instead of directly applying KD to the quantized student network from the beginning. Second, Co-studying (CS) phase tries to train a teacher to make it more-quantization-friendly and powerful than a fixed teacher. Finally, Tutoring (TU) phase transfers knowledge from the trained teacher to the student. This phase saves unnecessary training time and memory of the teacher network which tends to have already saturated in the co-studying phase. Fig. 11 shows the overall process of QKD. Self-studying (SS) phase gives a good starting point to alleviate the low representative power and the regularization effect of KD. Co-studying (CS) phase makes a teacher model adaptable to a student model and thereby powerful than the fixed teacher. In tutoring (TU)

phase, the teacher model transfers its adaptable and powerful knowledge to the student.

## IV. DEEP LEARNING PARADIGM

Knowledge distillation is flexibly applied to various deep learning paradigms since it was first based on the supervision of a teacher network. It is now expanding to the broad usage with weakly-supervised learning methodology and newly spotlighted unsupervised learning approach.

### 1) KD in Supervised Learning

Knowledge Distillation was originally introduced in the field of Supervised Learning [4], [6]. The vast majority of applications of Knowledge Distillation has been in this field to train a smaller compact network (student) under the supervision of a larger pre-trained network or an ensemble of models (teacher). In the context of deep learning, knowledge transfer has been successfully used to effectively compress the power of a larger capacity model (a teacher) to a smaller neural network (a student).
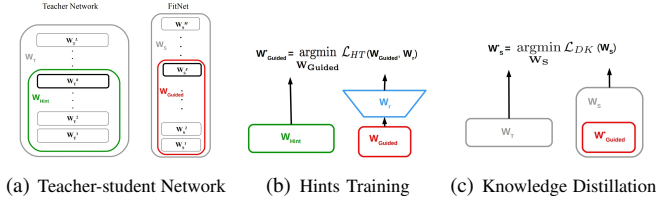


(a) Teacher-student Network  (b) Hints Training  (c) Knowledge Distillation

Fig. 12. Training a student network using hints [7].

---

**Algorithm 3:** FitNet Stage-Wise Training

**Input:** $\mathbf{W_S}, \mathbf{W_T}, g, h$

**Output:** $\mathbf{W_S^*}$

1 $\mathbf{W_{Hint}} \leftarrow \{\mathbf{W_T^1}, \ldots, \mathbf{W_T^h}\}$;
2 $\mathbf{W_{Guided}} \leftarrow \{\mathbf{W_S^1}, \ldots, \mathbf{W_S^g}\}$;
3 Initialize $\mathbf{W_r}$ to small random values;
4 $\mathbf{W_{Guided}^*} \leftarrow \arg\min_{\mathbf{W_{Guided}}} \mathcal{L}_{HT}(\mathbf{W_{Guided}}, \mathbf{W_r})$ ;
5 $\{\mathbf{W_S^1}, \ldots, \mathbf{W_S^g}\} \leftarrow \{\mathbf{W_{Guided}^{*1}}, \ldots, \mathbf{W_{Guided}^{*g}}\}$;
6 $\mathbf{W_S^*} \leftarrow \arg\min_{\mathbf{W_S}} \mathcal{L}_{KD}(\mathbf{W_S})$ ;

---

As an example, the work in Adriana Romero *et al.* [7] aimed to address the network compression problem by taking advantage of deep neural networks. It proposed an approach to train thin but deep neural networks, called FitNets, to compress wide and shallower (but still deep) networks. The method was extended to allow for thinner and deeper student models. In order to learn from the intermediate representations of teacher networks, FitNet made the student mimic the full feature maps of the teacher. However, such assumptions are too strict since the capacities of teacher and student may differ greatly. Fig. 12 shows how to train a student network using hints and Algorithm 3 shows the FitNet training algorithm. The algorithm receives as input the trained parameters $\mathbf{W_T}$ of a teacher, the randomly initialized parameters $\mathbf{W_S}$ of a FitNet, and two indices $h$ and $g$ corresponding to hint/guided layers, respectively. Let $\mathbf{W_{Hint}}$ be the teacher's parameters up to the hint layer $h$. Let $\mathbf{W_{Guided}}$ be the FitNet's parameters up to the guided layer $g$. Let $\mathbf{W_r}$ be the regressor's parameters. The first stage consists of pre-training the student network up to the guided layer, based on the prediction error of the teacher's hint layer (line 4). The second stage is a KD training of the whole network (line 6).

Sergey Zagoruyko *et al.* [8] proposed Attention Transfer (AT) to relax the assumption of FitNet. They transferred to the student network the attention maps that summarize the full activations.

### 2) KD in Weakly-Supervised Learning

Deep learning is data-hungry. Moreover, in supervised learning, labels corresponding to input data are required to train the network. However, it is very laborious to prepare for the label sets of training data. Weakly supervised learning [67] has been used to mitigate this problem. Yunchao Wei *et al.* [69] proposed *weakly supervised object detection* (WSOD) using Knowledge Distillation. Typically, fully-supervised object detection [70], [71] requires both labels for an object positional information (Anchor or Prior) and a class. In WSOD, only the class label is required for mining high-confidence region proposals with positive image-level annotations. The paper [69] utilized object segmentation knowledge to take benefit of WSOD.
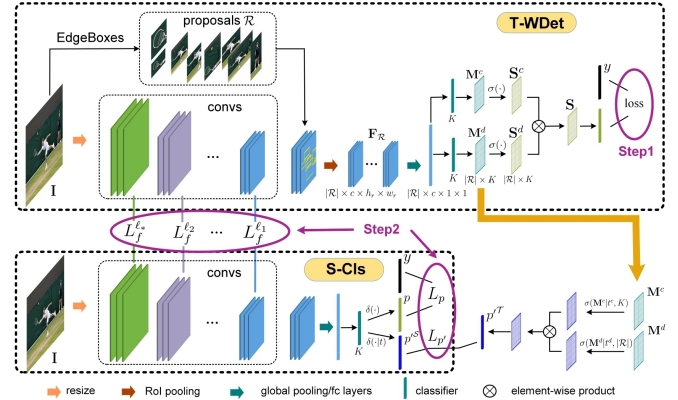


Fig. 13. An overall architecture [55]

Yongcheng Liu *et al.* [55] proposed a novel and efficient deep framework to boost multi-label classification by distilling knowledge from weakly-supervised detection task without bounding box annotations. This is an example of *cross-task knowledge distillation*. Specifically, given the image-level annotations, (1) a weakly-supervised detection (WSD) model is developed first, and then (2) an end-to-end multi-label image classification framework is constructed which is augmented by a knowledge distillation module that guides the classification model by the WSD model according to the class-level predictions for the whole image and the object-level visual features for object RoIs. The WSD model is the *teacher model* and the classification model is the *student model*. Fig. 13 shows the overall architecture of the network. The proposed framework works with two steps: (1) we first develop a WSD model as the teacher model (called T-WDet)

with only image-level annotations y; (2) then the knowledge in T-WDet is distilled into the MLIC student model (called S-Cls) via feature-level distillation from RoIs and prediction-level distillation from the whole image, where the former is conducted by optimizing the loss while the latter is conducted by optimizing the loss $L_p$ and $L'_p$.

*3) KD in Semi-Supervised Learning*

In the deep learning community, many studies agree that unsupervised learning is the future of deep learning. According to Yoshua Bengio and Yann LeCun, self-supervised learning – one of unsupervised learning approach - could lead to the creation of AI that's more humanlike in its reasoning. Recently, knowledge distillation is finding its ways in unsupervised learning [59]–[61].
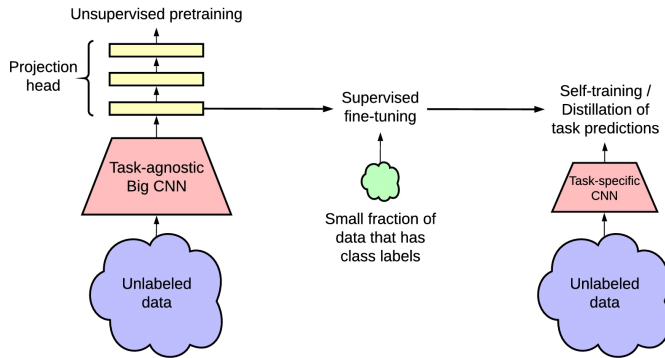


Fig. 14. **The proposed semi-supervised learning framework leverages unlabeled data in two ways:** (1) task-agnostic use in unsupervised pre-training, and (2) task-specific use in self-training / distillation [59].

Ting Chen *et al.* [59] showed that the paradigm of unsupervised pre-training followed by supervised fine-tuning is surprisingly effective for semi-supervised learning on ImageNet. The proposed semi-supervised learning algorithm can be summarized in three steps: unsupervised fine-tuning on a few labeled examples, and distillation with unlabeled examples for refining and transferring the task-specific knowledge. Fig. 14 shows the proposed framework.

## V. RELATED WORKS

In the perspective of machine teaching, there are a few of techniques related with knowledge distillation. Privileged information is the core concept of the smart teaching and this concept distinguishes it from knowledge distillation. David Lopez-Paz *et al.* [19] tried to unify knowledge distillation and privileged information approach and analyze them analytically.

### A. Learning Using Privileged Information (LUPI)

Vladimir Vapnik *et al.* [23] incorporated an "intelligent teacher" into machine learning. Their solution is to consider training data formed by a collection of triplets.

$$\{(x_1, x_1^*, y_1), \ldots, (x_n, x_n^*, y_n)\} \sim P^n(x, x^*, y), \quad (1)$$

where each $(x_i, y_i)$ is a feature-label pair, and the novel element $x_i^*$ is additional information about the example $(x_i, y_i)$

provided by an intelligent teacher, such as to support the learning process..

Even though, an additional information about the feature-label is provided by an intelligent teacher, the learning machine will not have an access to the teacher explanations $x_i^*$ at test time.

The framework of Learning Using Privileged Information (LUPI) [22], [23] studies how to leverage the additional information $x_i^*$ at training time to build a classifier for test time that outperforms those built on the regular features alone. LUPI has presented a new direction in knowledge transfer by modeling the transfer of prior knowledge as a Teacher-Student interaction process. Under LUPI, a Teacher model uses Privileged Information (PI) that is only available at training time to improve the sample complexity required to train a Student learner for a given task. At a high level, PI provides some similarity information between training samples from the original feature spaces, and the Teacher hypothesis serves as additional "explanations" of the hypothesis space.

Yunpeng Chen *et al.* [24] considered how to use PI to promote inherent diversity of a single CNN model such that the model can learn better representation and offer stronger generalization ability. To this end, they proposed a novel group orthogonal convolutional neural network (GoCNN) that learns untangled representations with in each layer by exploiting provided privileged information and enhances representation diversity effectively.

In this work, they proposed to exploit object segmentation annotations which are (partially) available in several public datasets as a privileged information for identifying the proper groups to give richer information. In addition, the background contents are usually independent on foreground objects within an image. Fig. 15 shows the architecture of the GoCNN. GoCNN is built upon a standard CNN architecture where the final convolution layer are explicitly divided into two groups: the foreground group (blue) which concentrates on learning the foreground feature and the background group (purple) which learns the background feature. The output features of these two groups are concatenated as a whole representation of the input image. In testing phase, parts within the gray shadow are removed and the "Concat" (green) operation is replaced by a "Pooling" operation making GoCNN back to a standard CNN.

There are some informative papers in the theoretical and applicable aspects of the LUPI framework. D. Pechyony *et al.* [26] gave a theoretical analysis of the LUPI framework and later M. Lapin *et al.* [27] showed that LUPI is a particular instance of importance weighting. The framework of LUPI also enjoyed multiple applications including ranking [28], computer vision [29], [30], clustering [31], metric learning [32], and Gaussian process classification [33].

### B. Generalized Distillation

Knowledge distillation [6] and privileged information [23] are two techniques that enable machines to learn from other machines.
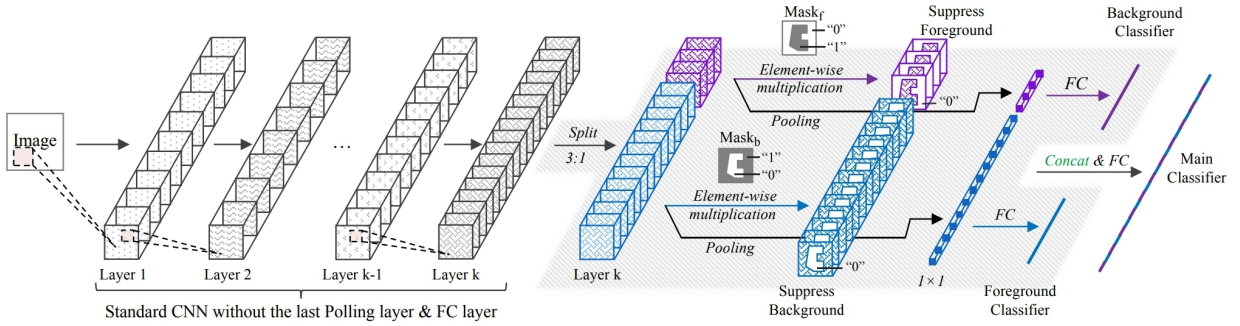
Fig. 15. The architecture of GoCNN [24].

There is an attempt to unify the two into generalized distillation [19] by casting distillation as a form of learning using privileged information, a learning setting in which additional per-instance information is available at training time but not test time. In short, generalized distillation is a framework to learn from multiple machines and data representations. David Lopez-Paz *et al.* [19] proposed the process of generalized distillation as follows and showed that generalized distillation reduces to knowledge distillation if $x_i^* = x_i$ for all $i$ with some constraints and it reduces to Vapnik's learning using privileged information if $x_i^*$ is a privileged description of $x_i$ with some constraints.

1) Learn teacher $f_t \in \mathcal{F}_t$ using the input-output pairs $\{(x_i^*, y_i)\}_{i=1}^n$ and (2).
2) Compute teacher soft labels $\{\sigma(f_t(x_i^*)/T)\}_{t=1}^n$, using temperature parameter $T > 0$.
3) Learn student $f_s \in \mathcal{F}_s$ using the input-output pairs $\{(x_i, y_i)\}_{i=1}^n$, $\{(x_i, s_i)\}_{i=1}^n$, (3) and imitation parameter $\lambda \in [0, 1]$.

$$f_t = \arg\min_{f \in \mathcal{F}_t} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, \sigma(f(x_i))) + \Omega(||f||) \tag{2}$$

$$f_s = \arg\min_{f \in \mathcal{F}_s} \frac{1}{n} \sum_{i=1}^n ((1-\lambda)\mathcal{L}(y_i, \sigma(f(x_i))) + \lambda\mathcal{L}(s_i, \sigma(f(x_i)))) \tag{3}$$

David Lopez-Paz *et al.* [19] also provided theoretical and causal insight about the inner workings of generalized distillation and showed some experiments to illustrate when the distillation of privileged information is effective, and when it is not.

## VI. FUTURE WORKS

### A. Explainable KD

In contrast to the empirical success of knowledge distillation, there is no satisfactory theoretical explanation of this phenomenon. For example, Seyed-Iman Mirzadeh *et al.* [17] showed empirically the effectiveness of introducing an intermediate network between student and teacher networks. They demonstrated in the experiment that TA with any intermediate size always improves the knowledge distillation performance. However, one might ask *"What the optimal TA size for the highest performance gain is?"*, *"If one TA improves the distillation results, why not also train this TA via another distilled TA?"*. They only showed feasibility with experiments and explained the results from empirical perspectives.

To our knowledge, David Lopez-Paz *et al.* [19] and Mary Phuong *et al.* [20] are the only works that examine distillation from a theoretical perspective. However, even the LUPI view conceptually falls short of explaining the effectiveness of distillation [19]. In particular, it concentrates on the aspect that the teacher's supervision to the student network is noise-free. Mary Phuong *et al.* [20] provides the first insights into the working mechanisms of distillation by studying the special case of linear and deep linear classifiers. Specifically it proves a generalization bound that establishes fast convergence of the expected risk of a distillation-trained linear classifier.

The mathematical principles underlying distillation's effectiveness have largely remained unexplored yet so this should be performed as a future research work.

### B. KD in Self-Supervised Image Representation Learning

In spite of the astonishing success of supervised learning, many deep learning researchers agree that unsupervised learning is the future of deep learning because the critical shortcoming of the supervised learning is the necessity of labeled big data. This can become a big huddle when we try to apply the deep learning approach to a new field where sufficient labeled data do not exist.

In machine learning, self-supervised learning has emerged as a paradigm to learn general data representations from

unlabeled examples and to fine-tune the model on labeled data. This has been particularly successful for natural language processing [62] and is an active research area for computer vision [63].

Recently, many studies pay attention to a variant of self-supervised learning in image understanding where the paradigm of unsupervised pre-training followed by supervised fine-tuning is surprisingly effective for semi-supervised learning on ImageNet [60]. Knowledge distillation is actively finding its way to this promising approach as an important component in its network pipeline. Ting Chen *et al.* [60] proposed the semi-supervised learning algorithm which can be summarized in three steps: unsupervised pre-tuning of a big ResNet model using SimCLRv2, supervised fine-tuning on a few labeled examples, and distillation with unlabeled examples for refining and transferring the task-specific knowledge.

## VII. CONCLUSION

The chronicled and comparative survey of the representative methods provide us holistic understanding and deep insight into knowledge distillation, which also motivate us to explore promising future directions.

The distillation-based approach to model compression has been proposed over a decade ago by C. Buciluǎ *et al.* [4] but was re-popularized by Hinton *et al.* [6], where additional intuition about why it works – due to the additional supervision and regularization of the higher entropy soft targets – was presented.

Knowledge distillation (KD) has proven to be a promising way to induce a small model that retains the accuracy of a large model but has the smaller computational complexity. It has been working by adding a distillation loss to the usual task loss so as to encourage the student network to mimic the teacher network's behavior. Nevertheless, a clear understanding of where valuable knowledge resides in a deep neural network is still lacking, and an optimal solution of how to capture the knowledge from a teacher network and transfer it to a student network remains an open question.

## REFERENCES

[1] Ian Goodfellow, Yoshua Bengio, Aaron Courville, Francis Bach. Deep Learning. MIT Press, 2016.
[2] Yu Cheng, Duo Wang, Pan Zhou. A Survey of Model Compression and Acceleration for Deep Neural Networks. IEEE Signal Processing Magazine, 2019.
[3] Zhengxia Zou, Zhenwei Shi, Jieping Ye. Object Detection in 20 Years: A Survey. arXiv preprint, 2019.
[4] C. Bucilua, R. Caruana, and A. Niculescu-Mizil. Model Compression. In Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '06, pages 535-541, New York, NY, USA, 2006. ACM.
[5] Ba J, and R. Caruana. 2014. Do Deep Nets Really Need to be Deep ?, In NIPS, 2654-2662
[6] G. Hinton, O. Vinyals, and J. Dean. 2015. Distilling the Knowledge in a Neural Network. In NIPS Deep Learning and Representation Learning Workshop.
[7] Adriana Romero, Nicolas Ballas. FitNets: Hints for thin deep nets, 2015
[8] Sergey Zagoruyko, Nikos Komodakis. Paying more attention to attention: Improving the performance of CNN via attention transfer, 2017
[9] Asit Mishra, Debbie Marr. Apprentice: Using KD techniques to improve low-precision network accuracy, 2017
[10] Antonio Polino, Razvan Pascanu. Model compression via distillation and quantization, 2018
[11] Jangho Kim, Yash Bhalgat. QKD: Quantization-aware knowledge distillation, 2019
[12] Guobin Chen, Wongun Choi. Learning efficient object detection models with knowledge distillation, 2017
[13] Quanquan Li, Shengying Jin. Mimicking very efficient network for object detection, 2017
[14] Wanwei Wang, Wei Hong. GAN-Knowledge distillation for one-stage object detection, 2019
[15] Shan You, Chang Xu, Chao Xu, Dacheng Tao. Learning from multiple teacher networks, 2017
[16] Ying Zhang, Tao Xiang. Deep Mutual Learning, 2017
[17] Seyed-Iman Mirzadeh, Mehrdad Farajtabar. Improved knowledge Distillation via Teacher Assistant, 2019
[18] Jiaxi Tang, Rakesh Shivanna, Zhe Zhao, Dong Lin, Anima Singh, Ed H. Chi, Sagar Jain. Understanding and Improving Knowledge Distillation, 2020
[19] David Lopez-Paz, Leon Bottou, Unifying distillation and privileged information. 2016
[20] Mary Phuong, Christoph H. Lampert. Towards understanding knowledge distillation, 2019
[21] Tao Wang, Li Yuan, Xiaopeng Zhang. Distilling object detectors with fine-grained feature imitation, 2019
[22] Vladimir Vapnik, Akshay Vashist. A new learning paradigm: Learning using privileged information. Neural Networks, 22(5):544-557, 2009
[23] Vladimir Vapnik, Rauf Izmailov. Learning using privileged information: Similarity control and knowledge transfer. JMLR, 16:2023-2049, 2015
[24] Yunpeng Chen, Xiaojie Jin, Shuicheng Yan, Training Group Orthogonal Neural Networks with Privileged Information. 2017
[25] Fengyi Tang, Cao Xiao, Fei Wang, Jiayu Zhou, Li-wei H. Lehman, Retraining Privileged Information for Multi-Task Learning. 2019
[26] Dmitry Pechyony, Vladimir Vapnik. On the theory of learning with privileged information. In NIPS, 2010.
[27] Maksim Lapin, Matthias Hein, Bernt Schiele. Learning using privileged information: SVM+ and weighted SVM. Neural Networks, 53:95-108, 2014.
[28] Viktoriia Sharmanska, Novi Quadrianto, Christoph H Lampert. Learning to rank using privileged information. In ICCV, 2013.
[29] Viktoriia Sharmanska, Novi Quadrianto, Christoph H Lampert. Learning to transfer privileged information. arXiv, 2014.
[30] David Lopez-Paz, Suvrit Sra, Alex Smola, Zoubin Ghahramani. Randomized nonlinear component analysis. In ICML, 2014.
[31] Jan Feyereisl, Uwe Aickelin. Privileged information for data clustering. Information Sciences, 194:4-23, 2012.
[32] Shereen Fouad, Peter Tino, Somak Raychaudhury, Petra Schneider. Incorporating privileged information through metric learning. Neural Networks and Learning Systems, 24(7):1086-1098, 2013.
[33] Daniel Hernandez-Lobato, Viktoriia Sharmanska, Kristian Kersting, Christoph H Lampert, and Novi Quadrianto. Mind the nuisance: Gaussian process classification using privileged noise. In NIPS, 2014.
[34] Byeong Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, Jin Young Choi. A Comprehensive Overhaul of Feature Distillation. In ICCV, 2019.
[35] A. Krizhevsky, I. Sutskever, and G. Hinton, "Imagenet classification with deep convolutional neural networks," in NIPS, 2012.
[36] Y. Gong, L. Liu, M. Yang, and L. D. Bourdev, "Compressing deep convolutional networks using vector quantization," CoRR, vol.abs/1412.6115, 2014.
[37] Y. W. Q. H. Jiaxiang Wu, Cong Leng and J. Cheng, "Quantized convolutional neural networks for mobile devices," in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.
[38] V. Vanhoucke, A. Senior, and M. Z. Mao, "Improving the speed of neural networks on cpus," in Deep Learning and Unsupervised Feature Learning Workshop, NIPS 2011, 2011.
[39] S. Gupta, A. Agrawal, K. Gopalakrishnan, and P. Narayanan, "Deep learning with limited numerical precision," in Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37, ser. ICML'15, 2015, pp. 1737–1746.
[40] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," International Conference on Learning Representations (ICLR), 2016.

[41] Y. Choi, M. El-Khamy, and J. Lee, "Towards the limit of network quantization," CoRR, vol. abs/1612.01543, 2016

[42] M. Courbariaux, Y. Bengio, and J. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," in Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada, 2015, pp. 3123–3131.

[43] M. Courbariaux and Y. Bengio, "Binarynet: Training deep neural networks with weights and activations constrained to +1 or -1," CoRR, vol.abs/1602.02830, 2016.

[44] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," in ECCV, 2016.

[45] P. Merolla, R. Appuswamy, J. V. Arthur, S. K. Esser, and D. S. Modha, "Deep neural networks are robust to weight binarization and other nonlinear distortions," CoRR, vol. abs/1606.01981, 2016.

[46] Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. PAMI, 2015

[47] Cai, Z., Vasconcelos, N.: Cascade R-CNN: High quality object detection and instance segmentation. PAMI, 2019.

[48] Lin, T.Y., Goyal, P., Girshick, R.B., He, K., Dollar, P.: Focal loss for dense object detection. In: ICCV, 2017.

[49] Zhou, X., Wang, D., Krahenbuhl, P.: Objects as points. arXiv:1904.07850, 2019.

[50] Tian, Z., Shen, C., Chen, H., He, T.: FCOS: Fully convolutional one-stage object detection. In: ICCV, 2019.

[51] Zhang, S., Chi, C., Yao, Y., Lei, Z., Li, S.Z.: Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection. arXiv:1912.02424, 2019.

[52] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End Object Detection with Transformer. arXiv:2005.12872v3, 2020.

[53] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, Antonio Torralba. Object Detectors Emerge in Deep Scene CNNs. arXiv:1412.6856v2, 2015.

[54] Byungseok Roh, Han0Cheol Cho, Myung-Ho Ju, and Soon Hyung Pyo. BABO: Background Activation Black-Out for Efficient Object Detection. arXiv:2002.01609v2, 2020.

[55] Yongcheng Liu, Lu Sheng, Jing Shao, Junjie Yan, Shiming Xiang, Chunhong Pan. Multi-Label Image Classification via Knowledge Distillation from Weakly-Supervised Detection. arXiv:1809.05884, 2019.

[56] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, Qun Liu. TinyBERT: Distilling BERT for Natural Language Understanding. arXiv:1909.10351, 2019.

[57] Victor SANH, Lysandre BEBUT, Julien CHAUMOND, Thomas WOLF. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. arXiv:1910.01108, 2020.

[58] Amjad Almahairi, Nicolas Ballas, Tim Cooijmans, Yin Zheng, Hugo Larochelle, Aaron Courville. Dynamic Capacity Networks. arXiv.1511.07838, 2016.

[59] Ting Chen, Simon Kornblith, Mohammad Norouzi, Geoffrey Hinton. A Simple framework for contrastive learning of visual representation. arXiv:2002.05709, 2020a.

[60] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, Geoffrey Hinton. Big Self-Supervised Models are Strong Semi-Supervised Learners. arXiv:2006.10029, 2020b.

[61] Jean-Bastien Grill, Florian Strub, Florent Altche,Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Remi Munos, Michal Valko. Bootstrap Your Own Latent: A New Approach to Self-Supervised Learning. arXiv:2006.07733, 2020.

[62] M.E.Peters, M. Neumann, M.Iyyer, M.Gardner, C.Clark, K.Lee, L.Zettlemoyer. Deep Contextualized Word Representations. In Proc. Of ACL, 2018.

[63] O.J.Henaff, A.Razavi, C.Doersch, S.M.A.Eslami, A. van den Oord. Data-efficient image recognition with contrastive predictive coding. arXiv, abs/1905.09272, 2019.

[64] O.Russakovsky, J.Deng, H.Su, J.Krause, S.Satheesh, S.Ma, Z.Huang, A.Karpathy, A.Khosla, M.Bernstein, A.C.Berg, L.Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. IJCV, vol.115, no. 3, pp. 211'252, 2015.

[65] Taehoon Kim, YoungJoon Yoo, Jihoon Yang. StatAssist & GradBoost: A Study on Optimal INT8 Quantization-aware Training from Scratch. arXiv:2006.09679v1, 2020.

[66] Muyang Li, Ji Lin, Yaoyao Ding,Zhijian Liu, Jun-Yan Zhu, and Song Han. Gan compression: Efficient architectures for interactive conditional gans. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020.

[67] Hakan Bilen, Andrea Vedaldi, "Weakly Supervised Deep Detection Networks". In: IEEE Computer Vision and Pattern Recognition, 2016.

[68] Yi Wei, Xinyu Pan,Hongwei Qin, and Junjie Yan. Quantization mimic: Towards very tiny CNN for object detection. CoRR, abs/1805.02152, 2018.

[69] Yunchao Wei, Zhiqiang Shen, Bowen Cheng, Honghui Shi, Jinjun Xiong, Jiashi Feng, Thomas Huang. TS2C: Tight Box Mining with Surrounding Segmentation Context for Weakly Supervised Object Detection. ECCV 2018, 2018.

[70] Shaoqing Ren, Kaiming He, Ross Girshick. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. arXiv:1506.01497, 2015.

[71] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. SSD: Single Shot Multibox Detector. In ECCV, 2016.

[72] Czarnecki, Osindero, Jaderberg, Swirszcz, Pascanu. Sobolev training for neural networks. In Advances in Neural Information Processing Systems, pp. 4278'4287, 2017.

[73] Srinivas, Fleuret. Knowledge transfer with jacobian match. arXiv:1803.00443, 2018.

[74] Junho Yim, Donggyu Joo, Jihoon Bae, Junmo Kim. A Gift from Knowledge Distillation: Fast Optimization, Network Minimization and Transfer Learning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4133-4141, 2017.

[75] Fahad Sarfraz, Elahe Arani and Bahram Zonooz. Knowledge Distillation Beyond Model Compression. arXiv:2007.01922v1, 2020.

[76] Guo, Qiushan, Xinjiang Wang, Yichao Wu, Zhipeng Yu, Ding Liang, Xiaolin Hu, and Ping Luo. "Online Knowledge Distillation via Collaborative Learning." In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11020-11029. 2020.

[77] Chen, Defang, Jian-Ping Mei, Can Wang, Yan Feng, and Chun Chen. "Online Knowledge Distillation with Diverse Peers." In AAAI, pp. 3430-3437. 2020.
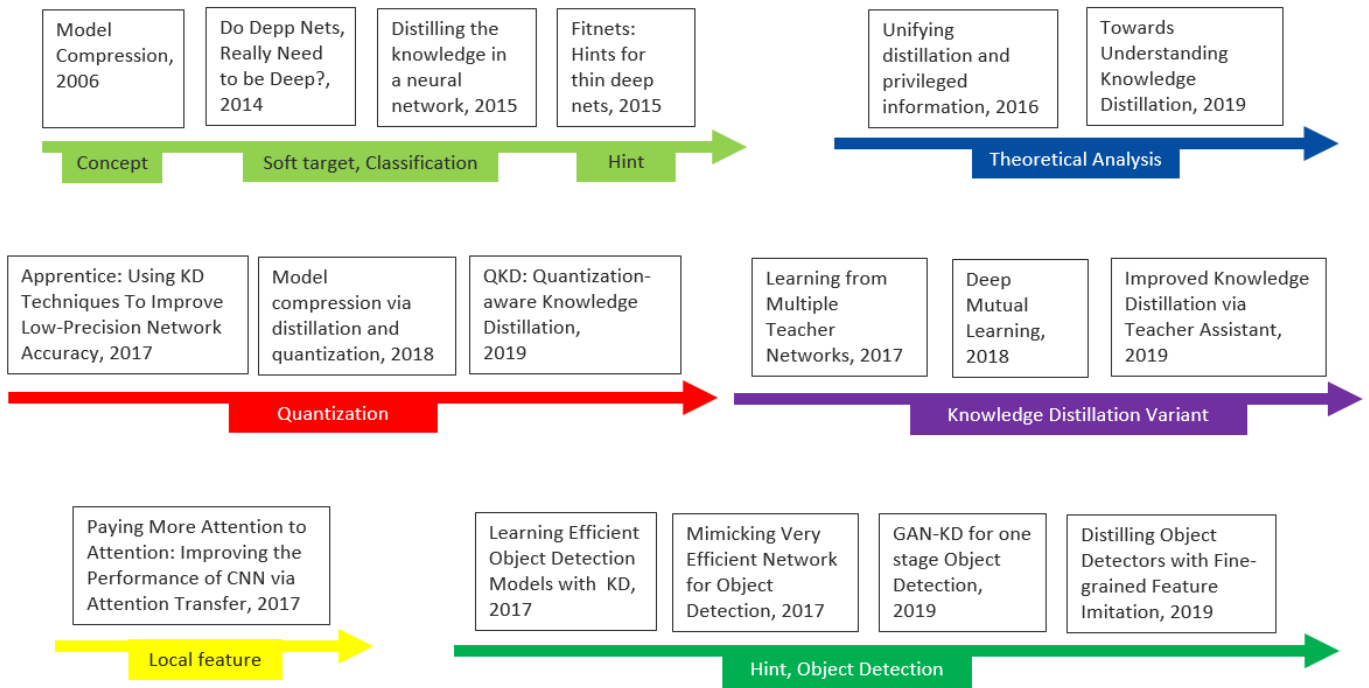
## ANNEX

### CONTENTS

Fig. 16. **Overall Knowledge Distillation (KD) Scope.** Overall Knowledge Distillation (KD) Scope. Overall scope of Knowledge Distillation (KD) and its related technical domains surveyed in this paper. Each arrow groups some related topics (represented as the title of the corresponding paper) chronically in that domain.
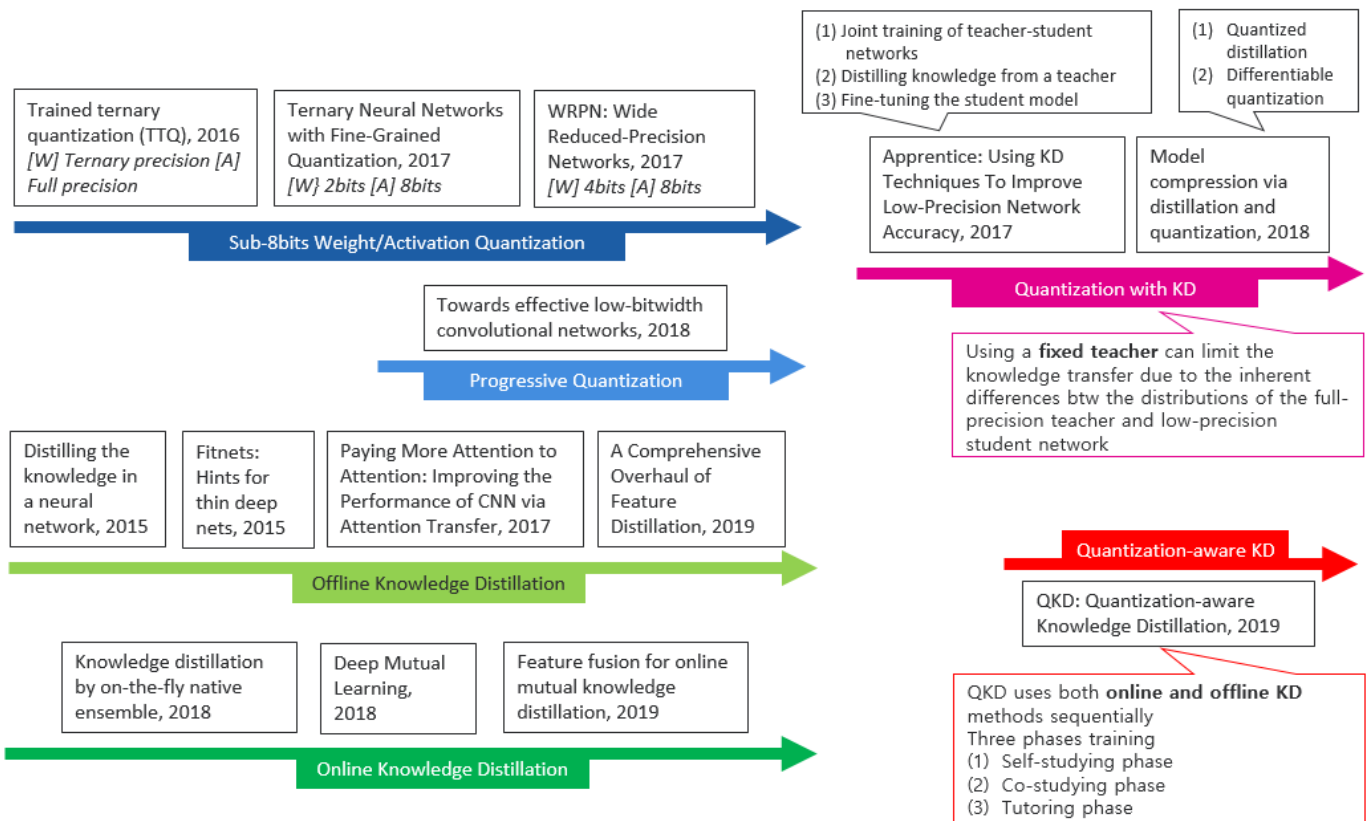


Fig. 17. **Knowledge Distillation Applied to Quantization.** Overallscope where knowledge distillation applied to quantization
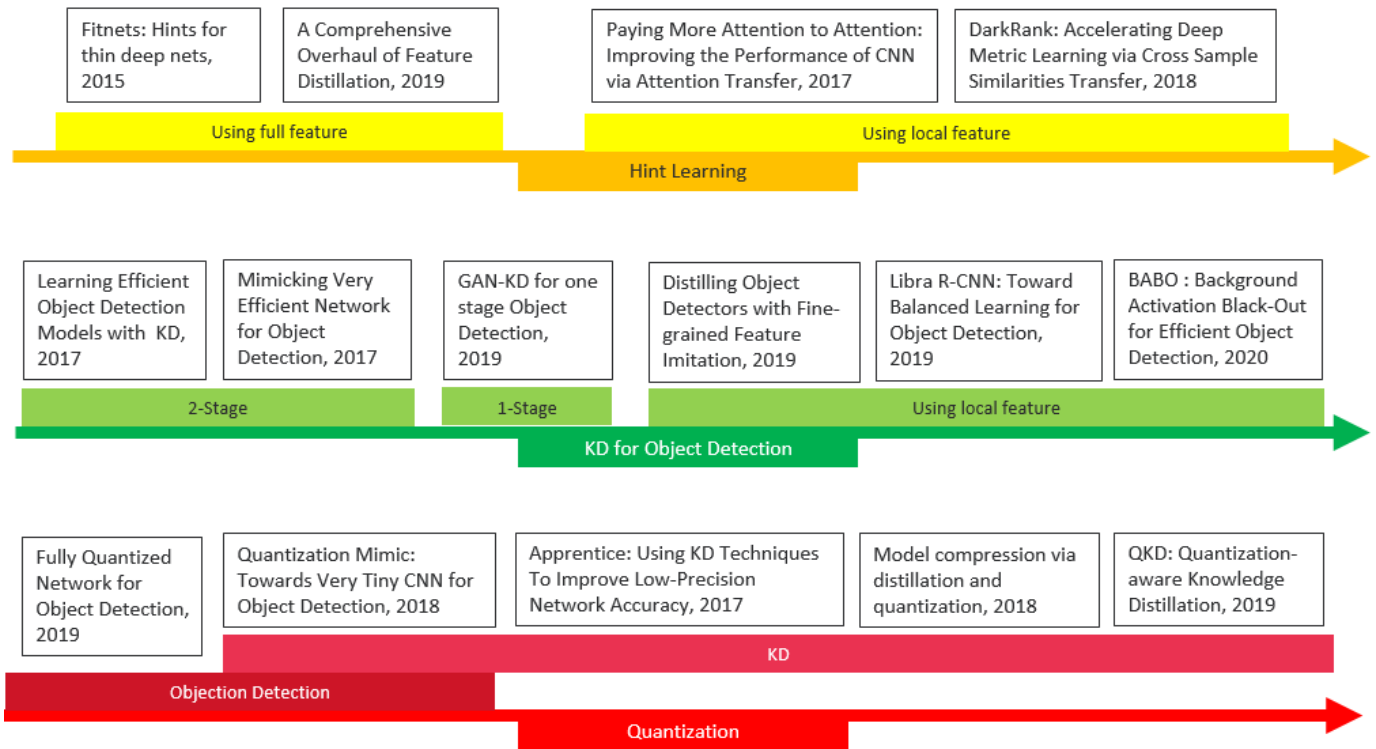
**Fitnets: Hints for thin deep nets, 2015**

**A Comprehensive Overhaul of Feature Distillation, 2019**

**Paying More Attention to Attention: Improving the Performance of CNN via Attention Transfer, 2017**

**DarkRank: Accelerating Deep Metric Learning via Cross Sample Similarities Transfer, 2018**

Using full feature — Using local feature

Hint Learning

**Learning Efficient Object Detection Models with KD, 2017**

**Mimicking Very Efficient Network for Object Detection, 2017**

**GAN-KD for one stage Object Detection, 2019**

**Distilling Object Detectors with Fine-grained Feature Imitation, 2019**

**Libra R-CNN: Toward Balanced Learning for Object Detection, 2019**

**BABO : Background Activation Black-Out for Efficient Object Detection, 2020**

2-Stage — 1-Stage — Using local feature

KD for Object Detection

**Fully Quantized Network for Object Detection, 2019**

**Quantization Mimic: Towards Very Tiny CNN for Object Detection, 2018**

**Apprentice: Using KD Techniques To Improve Low-Precision Network Accuracy, 2017**

**Model compression via distillation and quantization, 2018**

**QKD: Quantization-aware Knowledge Distillation, 2019**

KD

Objection Detection

Quantization

Fig. 18. **Knowledge Distillation Applied to Object Detection.** Overall scope where knowledge distillation is applied to object detection

Machines-teaching-machine paradigm

**Model Compression, 2006**

**Do Depp Nets, Really Need to be Deep?, 2014**

**Distilling the knowledge in a neural network, 2015**

**Fitnets: Hints for thin deep nets, 2015**

Knowledge Distillation (KD)

Unifies KD and PI into generalized distillation which is a framework to learn from multiple machines and data representation

**Unifying distillation and privileged information, 2016**

Generalized Distillation

**A New Learning Paradigm: Learning Using Privileged Information , 2009**

**On the Theory of Learning with Privileged Information, 2010**

PI is an additional information about the {feature, label} provided by an intelligent teacher

Privileged Information (PI)

Designing a fully data-driven automated TA selection ?

**"Why and when introducing a TA improves KD?"**
The student network performance degrades when the gap btw student and teacher is large
※ Theoretical analysis using VC theory

An optimization framework to obtain good teaching strategy.
The teacher model leverages the feedback from the student model to optimize its own teaching strategies

**Intelligent Tutoring System, 1985**

**Machine Teaching for Bayesian Learners in the Exponential Family, 2013**

**Learning to Teach, 2018**

Machine Teaching

**Learning from Multiple Teacher Networks, 2017**

**Deep Mutual Learning, 2018**

**Improved Knowledge Distillation via Teacher Assistant, 2019**

Knowledge Distillation Variant

Fig. 19. **Knowledge Distillation Variants, Learning Using Privileged Information (LUPI).** Overall scope which represents the relationship among knowledge distillation, learning with privileged information and generalized distillation