

Learning from Very Few Samples: A Survey

Jiang Lu, Pinghua Gong, Jieping Ye, *Fellow, IEEE*, Jianwei Zhang, *Member, IEEE*,
and Changshui Zhang, *Fellow, IEEE*

Abstract—Few sample learning (FSL) is significant and challenging in the field of machine learning. The capability of learning and generalizing from very few samples successfully is a noticeable demarcation separating artificial intelligence and human intelligence since humans can readily establish their cognition to novelty from just a single or a handful of examples whereas machine learning algorithms typically entail hundreds or thousands of supervised samples to guarantee generalization ability. Despite the long history dated back to the early 2000s and the widespread attention in recent years with booming deep learning technologies, little surveys or reviews for FSL are available until now. In this context, we extensively review 300+ papers of FSL spanning from the 2000s to 2019 and provide a timely and comprehensive survey for FSL. In this survey, we review the evolution history as well as the current progress on FSL, categorize FSL approaches into the generative model based and discriminative model based kinds in principle, and emphasize particularly on the meta learning based FSL approaches. We also summarize several recently emerging extensional topics of FSL and review the latest advances on these topics. Furthermore, we highlight the important FSL applications covering many research hotspots in computer vision, natural language processing, audio and speech, reinforcement learning and robotic, data analysis, etc. Finally, we conclude the survey with a discussion on promising trends in the hope of providing guidance and insights to follow-up researches.

Index Terms—few sample learning, learn to learn, survey, few-shot learning, meta learning

1 INTRODUCTION

ONE impressive hallmark of human intelligence is the ability to rapidly establish cognition to novel concepts from just a single or a handful of examples. Many cognitive and psychological evidences [1], [2], [3] have shown that humans can recognize visual objects through very few images [4] and even children can remember a novel word by a single encounter [5], [6]. Although exactly what support the human capability of learning and generalizing from very few samples remains a profound mystery, some neurobiological works [7], [8], [9] have argued that the prominent human learning ability benefits from prefrontal cortex (PFC) and working memory in human brain, especially the interaction between PFC-specific neurobiological mechanism and previous experience stored in the brain. By contrast, most cutting-edge machine learning algorithms are data-hungry, especially the most widely known deep learning [10] that has pushed artificial intelligence to a new climax. As an important milestone in the development of machine learning, deep learning has scored remarkable achievement in a broad spectrum of research fields including vision [11], [12], [13], language [14], [15], speech [16], game [17], demography [18], medicine [19], phytopathology [20] and zoology [21], etc. Generally, the successes of deep learning can be owned to three key factors: powerful computing resources (e.g., GPU), sophisticated

neural networks (e.g., CNN [11], LSTM [22]) and large-scale datasets (e.g., ImageNet [23], Pascal-VOC [24]). However, many realistic application scenarios, such as in the field of medicine, military and finance, do not allow us access sufficient labeled training samples, due to some factors including privacy, security or high labeling costs for data, etc. Thus, it becomes an eagerly-awaited blueprint for almost all machine learning researchers to enable learning systems to efficiently learn and generalize from very few samples.

From a high-level perspective, the theoretical and practical significance of studying few sample learning (FSL) mainly comes from three aspects. First, the FSL approach is expected not to rely on large-scale training samples, thus eschewing the prohibitive costs on data preparation in some specific applications. Second, FSL can shrink the gap between human intelligence and artificial intelligence, being a necessary trip to develop universal AI [25]. Third, FSL can achieve a low-cost and quick model deployment for one emerging task for which just a few samples are temporarily available, beneficial to shed light on the potential laws earlier in the task.

Despite these encouraging virtues, the research of FSL progresses more slowly in the past decades than that of large sample learning due to its intrinsic difficulty. Clearly, we illustrate this difficulty from the optimization viewpoint. Consider a general machine learning problem, which is described by a prepared supervised training set $\mathcal{D}_t = \{(x_i, y_i)\}_{i=1}^n$ with $x \in \mathcal{X}$, $y \in \mathcal{Y}$ drawn from the joint distribution $P_{\mathcal{X} \times \mathcal{Y}}$. The goal of the learning algorithm is to produce a mapping function $f \in \mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$ such that the expected error $\mathcal{E}_{\text{ex}} = \mathbb{E}_{(x,y) \sim P_{\mathcal{X} \times \mathcal{Y}}} L(f(x), y)$ is minimized, where $L(f(x), y)$ denotes the loss that compares the prediction $f(x)$ to its supervision target y . In fact, the joint distribution $P_{\mathcal{X} \times \mathcal{Y}}$ is unknown, and thus the learning algorithms are intended to minimize the empirical error $\mathcal{E}_{\text{em}} = \mathbb{E}_{(x,y) \sim \mathcal{D}_t} L(f(x), y)$. In this context, a typical problem is that if the function space \mathcal{F} from which the learning algorithm selects f is too large,

- *J. Lu and C. Zhang are with the Institute for Artificial Intelligence, Tsinghua University (THUAI), the State Key Laboratory of Intelligence Technologies and Systems, the Beijing National Research Center for Information Science and Technologies (BNRist), the Department of Automation, Tsinghua University, Beijing 100084, China. E-mail: lu.j13@tsinghua.org.cn; zcs@mail.tsinghua.edu.cn.*
- *P. Gong is with the Didi Research Institute, Didi Chuxing, Beijing 100085, China. E-mail: gongpinghua@didichuxing.com.*
- *J. Ye is with Didi AI Labs, Beijing 100085, China and University of Michigan, Ann Arbor. E-mail: jye@umich.edu.*
- *J. Zhang is with the Faculty of Mathematics Computer Science and Natural Sciences, TAMS Group, University of Hamburg, Hamburg 20146, Germany. E-mail: zhang@informatik.uni-hamburg.de.*

the generalization error $\mathcal{E} = |\mathcal{E}_{\text{ex}} - \mathcal{E}_{\text{em}}|$ would become big and thereby overfitting may arise easily. We can re-look the problem from the following perspective

$$\min_f \mathcal{E}_{\text{em}}, \quad \text{s.t. } f(x_i) = y_i, \forall (x_i, y_i) \in \mathcal{D}_t. \quad (1)$$

If \mathcal{D}_t contains more supervised samples, there will be more constraints on f , which implies the space of function f will be smaller, then it will bring a good generalization. Conversely, a scarce supervised training set would naturally lead to a poor generalization performance. Essentially, the constraint formed by each supervised sample can be regarded as a regularization on the function f , which is able to compress the redundant optional space of function f and thereby reduce its generalization error. Thus, it can be concluded that if one learning algorithm deal with one FSL task just by the vanilla learning techniques without any sophisticated learning strategies or specific network design, the learning algorithm would be faced with the serious overfitting.

Few sample learning (FSL), also known as small or one sample learning, few-shot or one-shot learning, can date back to the early 2000s. Despite the nearly 20 years of research history and its importance at the level of theory and application, few related surveys or reviews are available until now. In this article, we extensively investigate almost all FSL-related scientific papers spanning from the 2000s to 2019 to elaborate a systematic FSL survey. We must emphasize that the FSL discussed here is orthogonal to zero-shot learning (ZSL) [26], which is another hot topic for machine learning. The setting of ZSL entails concept-specific side information to support the cross-concept knowledge transfer, varying greatly from that of FSL. To our best knowledge, there only have two FSL-related preprinted surveys [27], [28] until now. Compared with them, the novelties and contributions of this survey mainly come from five major aspects:

(1) We give a more comprehensive and timely review which encompasses 300+ FSL-related papers spanning from the 2000s to 2019, covering all FSL approaches from the very earliest Congealing model [29] to the latest meta learning approaches. The exhaustive exposition is conducive to the grasp of the whole development process of FSL as well as the construction of the complete knowledge hierarchy to FSL.

(2) We provide an understandable hierarchical taxonomy that categorizes existing FSL approaches into the generative model based approaches and discriminative model based approaches in light of their modeling principles to FSL problems. Within each class, we further conduct a more detailed categorization according to the generalizable properties.

(3) We put emphasis on current mainstream FSL approaches, i.e., the meta learning based FSL approaches, and categorize them into five major classes in light of what they hope to learn to learn via meta learning strategy, including Learn-to-Measure, Learn-to-Finetune, Learn-to-Parameterize, Learn-to-Adjust and Learn-to-Remember. Moreover, the underlying development relationship between various meta learning based FSL approaches is revealed in this survey.

(4) We conclude several extensional research topics beyond vanilla FSL that are emerging lately and review the latest advances towards these topics. These topics include Semi-supervised FSL, Unsupervised FSL, Cross-domain FSL, Generalized FSL and Multimodal FSL, which are challenging whilst endowing prominent practical significance to the

solution for many realistic machine learning problems. These extensional topics were rarely covered by previous reviews.

(5) We extensively summarize existing FSL applications in various fields including computer vision, natural language processing, audio and speech, reinforcement learning and robotic, data analysis, etc, and current FSL performance on benchmarks, aiming to provide a handbook for follow-up researches, which were not studied by previous reviews.

The rest of this paper is organized as follows. In Section 2, we give a general overview including the evolution history of FSL, the notations and definition we will use later, and the proposed taxonomy for existing FSL approaches. The generative model based approaches and the discriminative model based approaches are discussed in detail in Section 3 and Section 4, respectively. Then, several emerging extensional topics for FSL are summarized in Section 5. In Section 6, we extensively investigate the FSL applications in various fields and the benchmark performance of FSL. In Section 8, we conclude this survey with a discussion on future directions.

2 OVERVIEW

In this section, we first briefly review the FSL evolution history in Section 2.1. Then, some notations and definitions are introduced in Section 2.2. Finally, we provide a high-level taxonomy for existing FSL approaches in Section 2.3.

2.1 Evolution History

The general regime of machine learning is to make predictions on the future data using the statistical models that are learned on previously prepared training samples. In most cases, the generalization ability of the models is guaranteed by a sufficient quantity of training samples. In many realistic applications, nevertheless, we might be allowed to access only very few training data for novel concepts, in the limit, just one example per concept. For instance, we may need to recognize several kinds of uncommon animals whereas only several annotated pictures are at hand due to their rarity. Similarly, we may be required to authenticate the identity for some new users based on mobile sensor information given a handful of historical usage records from them. The problem of learning from very few examples firstly attracted the attention of E. G. Miller *et al.* in 2000 [29], who postulated a shared density on digit transforms and proposed a Congealing algorithm to bring test digit image into correspondence with class-specific congealed digit image. Thereafter, more and more efforts were devoted to FSL research.

The development process of FSL research can be roughly divided into two periods, non-deep period (from 2000 to 2015) and deep period (from 2015 to now), as depicted in Fig. 1. The watershed separating them is the first combination of deep learning techniques and FSL problems introduced by G. Koch *et al.* in 2015 [30]. Before that, all solutions proposed for FSL problems are based on non-deep learning methodologies or techniques. In particular, most of the famous early FSL approaches in non-deep period are based on the generative model. They seek to estimate the joint distribution $P(\mathcal{X}, \mathcal{Y})$ or the conditional distribution $P(\mathcal{X}|\mathcal{Y})$ given a supervision (e.g., a class), albeit on very few observed training samples and then make predictions for test samples

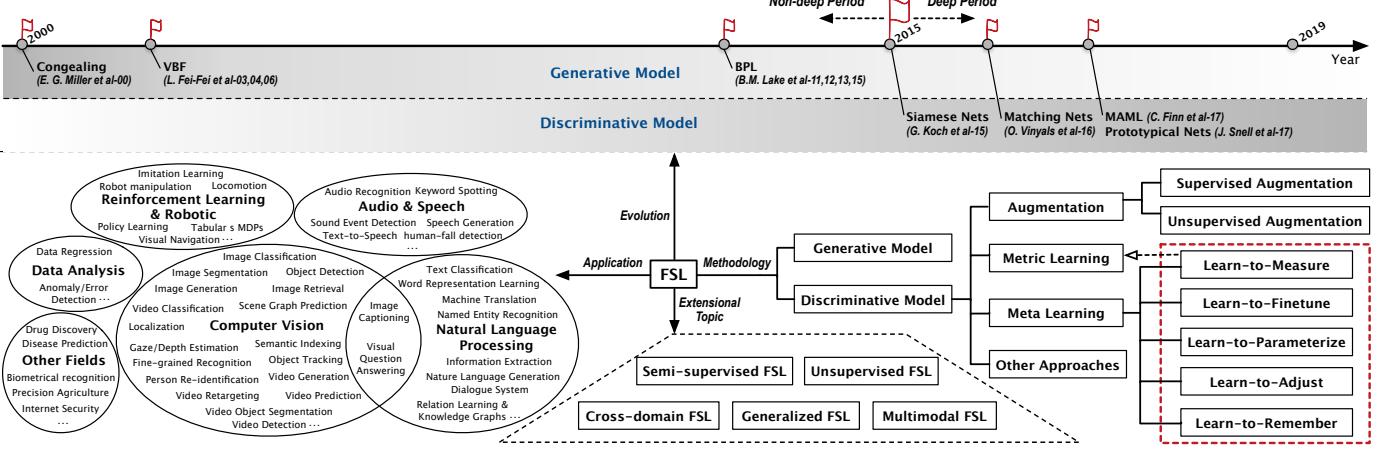


Fig. 1. Outline of our survey. The main contents include the evolution history, methodology, extensional topics, and applications of FSL.

from the point of Bayesian decision. Several milestones among these generative model based FSL approaches in non-deep period include Congealing algorithm by E. G. Miller *et al.* [29], Variational Bayesian framework (VBF) by L. Fei-Fei *et al.* [31], [32], [33], and Bayesian Program Learning (BPL) by B. M. Lake *et al.* [34], [35], [36], [37]. Congealing algorithm [29] is the earliest founder for studying how to learn from very few samples, while VBF [31] is the first work to articulate the term of “one-shot learning”. Comparably, BPL [37] reaches a human-level one-shot character classification performance by capitalizing on the human abilities of compositionality, causality and imagination in the cognition of novel concepts. In this non-deep period, there also have several discriminative model based FSL approaches [38], [39], [40], [41], [42], [43], though they were not the mainstream at this period. Opposite to generative model, discriminative model based FSL approaches pursue a conditional distribution $P(\mathcal{Y}|\mathcal{X})$ which can directly predict a probability given one observed sample. Despite the efforts above, the FSL research in the non-deep period still evolves very slowly.

With deep learning booming, especially the great success achieved by CNNs on visual tasks [11], [12], [13], many FSL researchers began to shift their sights from non-deep models to deep models. In 2015, G. Koch *et al.* [30] took the lead in incorporating deep learning into the solution for FSL issues by proposing a Siamese CNN to learn a class-irrelevant similarity metric on pairwise samples, which marks the beginning of a new era for FSL, i.e., the deep period. After that, the subsequent FSL approaches made full use of the advantages of deep neural networks in feature representation and end-to-end model optimization to address FSL problems from different angles including data augmentation [44], metric learning [45] and meta learning [46], etc, pushing FSL researches into a new period of rapid development. Although a few generative model based approaches, such as Neural Statistician [47] and Sequential Generative Model [48], were proposed in this deep period, discriminative model based FSL approaches dominate the evolution of FSL study. Especially, a large number of meta learning based FSL approaches have been springing up in recent years, such as Matching Nets by O. Vinyals *et al.* [49], MAML by C. Finn *et al.* [50], Meta-Learner LSTM by S. Ravi and H.

Larochelle [51], MANN by A. Santoro *et al.* [52], MetaNet by T. Munkhdalai and H. Yu [53], Prototypical Nets by J. Snell *et al.* [54], Relation Net by F. Sung *et al.* [55] and LGM-Nets by H. Li *et al.* [56], etc. Noticeably, meta learning strategies become the prevailing ideology for FSL. In this period, furthermore, these advanced FSL approaches have been directly applied to or improved to tackle various applications in computer vision, natural language processing, audio and speech, data analysis, robotics, etc. Meanwhile, more and more challenging extensional topics relating to FSL, such as Semi-supervised FSL, Unsupervised FSL, Cross-domain FSL, Generalized FSL and Multimodal FSL, have been unearthed.

In brief, the evolution history of FSL witnessed a transition from non-deep period to deep period, an alternation of mainstream approaches between generative model and discriminative model, and a resurgence of the classical meta learning idea. Today, FSL related works frequently appear in many top venues, most notably in machine learning or their applications, attracting wide attention of the machine learning community.

2.2 Notations and Definitions

Formally, we use x to represent input data, y to represent supervision target, \mathcal{X} and \mathcal{Y} to denote the space of input data and supervision target, respectively. An FSL task T is described by a T -specific dataset $D_T = \{D_{\text{trn}}, D_{\text{tst}}\}$ with $D_{\text{trn}} = \{(x_i, y_i)\}_{i=1}^{N_{\text{trn}}}$ and $D_{\text{tst}} = \{x_j\}_{j=1}^{N_{\text{tst}}}, x_i, x_j \in \mathcal{X}_T \subset \mathcal{X}, y_i \in \mathcal{Y}_T \subset \mathcal{Y}$. The samples x_i, x_j for task T come from one specific domain $\mathcal{D}_T = \{\mathcal{X}_T, P(\mathcal{X}_T)\}$ consisting of a data space \mathcal{X}_T and a marginal probability distribution $P(\mathcal{X}_T)$. Usually, there are C task classes and only K (very small, 1, 5, for example) samples per class in D_{trn} , that is, $N_{\text{trn}} = CK$, then T is also called as C -way K -shot task. The goal is to produce a target predictive function $f \in \mathcal{F} : \mathcal{X} \rightarrow \mathcal{Y}$ which can make predictions for test samples in D_{tst} . Based on our analysis in Section 1, it is hard to build a high-quality f just with the scarce D_{trn} . In most cases, therefore, researchers are allowed to leverage one supervised auxiliary dataset $D_A = \{(x_i^a, y_i^a)\}_{i=1}^{N_{\text{aux}}}, x_i^a \in \mathcal{X}_A \subset \mathcal{X}, y_i^a \in \mathcal{Y}_A \subset \mathcal{Y}$, that includes sufficient samples and classes ($N_a \gg N_{\text{trn}}, |\mathcal{Y}_A| \gg |\mathcal{Y}_T|$) collected based on previously seen concepts. It needs to be noted that D_A does not contain data belonging to classes in

T , that is, $\mathcal{Y}_T \cap \mathcal{Y}_A = \emptyset$, and the data of D_T and those of D_A come from the same domain, that is, $\mathcal{D}_T = \mathcal{D}_A$, $\mathcal{X}_T = \mathcal{X}_A$ and $P(\mathcal{X}_T) = P(\mathcal{X}_A)$, where $\mathcal{D}_A = \{\mathcal{X}_A, P(\mathcal{X}_A)\}$. The setting is solid and reasonable since D_A is easy to be acquired from many historical, offline or publicly well-labeled data that are relevant to task T , especially in today's big data era.

On these basis, we give a unified definition of FSL.

Definition 1 (Few Sample Learning) Given a task T described by a T -specific dataset D_T with only a few supervised information available, and a T -irrelevant auxiliary dataset D_A (if any), *few sample learning* aims to build a function f for task T that maps its inputs to targets using the very few supervision information in D_T and the knowledge in D_A .

The term of T -irrelevant in above definition implies the targets in D_T and D_A are orthogonal, that is, $\mathcal{Y}_T \cap \mathcal{Y}_A = \emptyset$. If D_A covers the classes in T , i.e., $\mathcal{Y}_T \cap \mathcal{Y}_A = \mathcal{Y}_T$, the FSL problem will collapse to a traditional large sample learning problem. In particular, if $|\mathcal{Y}_T| = 2$, T is a binary FSL task, and if $|\mathcal{Y}_T| > 2$, then we call T is a multiclass FSL task. Besides, we conclude several important extensional topics of FSL in light of the above notations and definition.

Semi-supervised FSL. In addition to the CK supervised samples, D_{trn} also contains some unlabeled training samples.

Unsupervised FSL. D_A is fully unsupervised despite it contains sufficient samples from non-task classes.

Cross-domain FSL. The samples in D_T and D_A come from two different data domains, that is, $\mathcal{D}_T \neq \mathcal{D}_A$.

Generalized FSL. Function f is required to make inference on united label space $\mathcal{Y}_T \cup \mathcal{Y}_A$ rather than single \mathcal{Y}_T .

Multimodal FSL. It has two cases, multimodal matching and multimodal fusion. In the former case, the target y_i in a labeled sample pair (x_i, y_i) is not the simple class label, but one data in another modality different from the modality of input x_i . In the latter case, the additional information I_i for x_i belonging to other modality is provided.

We give the detailed problem description and literature review for the above five extensional topics in Section 5.

2.3 Taxonomy

As shown in Fig. 1, we organize FSL approaches into two major categories, i.e., generative model based approaches and discriminative model based approaches in light of the modeling principles to FSL problems. For one test sample x_j , all FSL solutions pursue the following statistical model which can predict the posterior probability of class given x_j

$$\hat{y}_j = \arg \max_{y \in \mathcal{Y}_T} p(y|x_j), \quad (2)$$

where \hat{y}_j denotes the predicted target by this model. Discriminative model based FSL approaches aim to directly model the posterior probability $p(y|x)$, which takes x as the input of discriminative model and outputs one probability distribution of x belonging to C task classes. By contrast, generative model based approaches tackle it using Bayesian decision $p(y|x) = p(x|y)p(y)/p(x, y)$. Thus the maximization of posterior probability in Eq. (2) becomes to

$$\hat{y}_j = \arg \max_{y \in \mathcal{Y}_T} p(x_j|y)p(y), \quad (3)$$

where $p(y)$ is the prior distribution of target class, $p(x_j|y)$ is the conditional distribution of data given class y . In most

cases, $p(y)$ is assumed to be a uniform distribution among classes or computed as the frequency ratio of data in different classes. Consequently, the core aim of generative model based FSL approaches is to compare $p(x|y)$, $y \in \mathcal{Y}_T$.

For the category of generative model, researchers bridge the connection between x and y using some intermediate latent variables such that the conditional distribution $p(x|y)$ can be computed mathematically. Most of FSL approaches in this category require some necessary assumptions on the distribution of the latent variables. We will briefly discuss the generative model based approaches in Section 3.

For the category of discriminative model, three main-streams are summarized, which include augmentation, metric learning and meta learning. The augmentation approaches are further divided into supervised augmentation and unsupervised augmentation according to whether extra supervision information (e.g., attribute annotation, word embedding, etc) has been used. As the most popular treatment to FSL problems recently, the meta learning based approaches involve various views to reach the goal of learn-to-learn. We divide the existing meta learning based FSL approaches into five major genres in light of what is hoped to be meta-learned behind the meta strategy, Learn-to-Measure, Learn-to-Finetune, Learn-to-Parameterize, Learn-to-Adjust and Learn-to-Remember. In a broad sense, Learn-to-Measure approaches fall into the scope of metric learning since they all pursue a metric space rendering homogeneous samples close and inhomogeneous samples far apart. Even so, the important basis by which we assign Learn-to-Measure approaches into meta learning is the use of meta learning strategy. Also, there exist other niche directions to tackle FSL problems. We will review the discriminative model based FSL approaches in Section 4.

3 GENERATIVE MODEL BASED APPROACHES

As mentioned in Section 2.3, the generative model based FSL approaches seek to model the posterior probability $p(x|y)$. In most cases, however, the probabilistic relationship between data x and target y is not straightforward. For instance, in few-shot image classification, x denotes one image and y denotes its class label, and the mathematical connection between them can not be described directly. A feasible strategy to bridge the connection between x and y is to introduce an intermediate latent variable \mathbf{z} as follows:

$$p(x|y) = \int_{\mathbf{z}} p(x, \mathbf{z}|y) d\mathbf{z} = \int_{\mathbf{z}} p(\mathbf{z}|y)p(x|\mathbf{z}, y) d\mathbf{z}. \quad (4)$$

Almost all generative model based FSL approaches follow this high-level strategy, even if they differ in the specific form of \mathbf{z} . Several classic forms of \mathbf{z} are summarized as follows.

- Transformation

As the first work that attempts to learn from one sample, Congealing [29] algorithm assumes there exists one latent image for each digit class and all observed images belonging to this class are produced from the latent image through some underlying transformations \mathbf{z}_{tran} . Moreover, the density over transformations are supposed to be shared across different

TABLE 1
Summary of different generative model based FSL approaches

Approaches	Latent Variable	Task Type	Experimental Dataset	Remark
Congealing [29]	Transformation \mathbf{z}_{tran}	Multi-class image classification	NIST Special Database 19 [57]	the founder of FSL/only applicable to simple digit or letter character grayscale images
VBF [31], [32], [33]	Parameters \mathbf{z}_{para}	Binary image classification	Caltech 4 Data Set [31], [58], Caltech 101 Data Set [32], [33]	the first work to propose “one-shot learning”/ hard to adapt to multi-class tasks
HB [59]	Superclass \mathbf{z}_{sup}	Binary image classification	MNIST [60], MSR Cambridge dataset [59]	relies on the underlying hierarchical inter-class relationship/hard to adapt to multi-class tasks
BPL [34], [35], [36], [37]	Programs \mathbf{z}_{prog}	Multi-class image classification, Image generation	Omniglot [37]	requires the dynamic stroke information and the production rules of image objects
Chopping [61]	Splits \mathbf{z}_{spl}	Binary image classification	COIL-100 database [62], LATEX symbols [61]	like a probabilistic “ensemble” method/hard to adapt to multi-class tasks
CPM [63]	Reconstruction \mathbf{z}_{rec}	Multi-class image classification	MNIST [60], USPS [64]	only applicable to simple digit or letter character grayscale images/does not need the auxiliary set D_A
Neural Statistician [65]	Statistics \mathbf{z}_{stat}	Multi-class image classification, Image generation	MNIST [60], Omniglot [37], Youtube Faces database [66]	an extension of a variational autoencoder/ contains some deep neural networks

classes, which implies that the transformation probability is independent of class. Thus, Eq. (4) can be written into

$$p(x|y) = \int_{\mathbf{z}_{\text{tran}}} p(\mathbf{z}_{\text{tran}}) p(x|\mathbf{z}_{\text{tran}}, y) d\mathbf{z}_{\text{tran}}. \quad (5)$$

Importantly, $p(\mathbf{z}_{\text{tran}})$ can be learned on auxiliary set D_A . We must emphasize that Congealing algorithm is only applicable to the simple digit or letter character grayscale images since it is unrealistic to model such class-shared transformation mathematically for other natural RGB images.

- *Parameters*

VBF [31], [32], [33] measures the probability that an object exists in one RGB image using probabilistic models. The probabilistic models involve many parameters \mathbf{z}_{para} that needs to be learned. Thus, VBF defines the $p(\mathbf{z}_{\text{para}}|y)$ using a so-called constellation model and utilizes variational methods to estimate \mathbf{z}_{para} on auxiliary set D_A .

- *Superclass*

In [59], a hierarchical Bayesian (HB) model was developed by introducing the superclass relationship over classes. Its key insight is that the classes under the same superclass inherit the same similarity metric. By the superclass variable \mathbf{z}_{sup} , Eq. (4) can be turned into

$$p(x|y) = \sum_{\mathbf{z}_{\text{sup}}} p(\mathbf{z}_{\text{sup}}^y) p(x|\mathbf{z}_{\text{sup}}^y), \quad (6)$$

where $p(\mathbf{z}_{\text{sup}}^y) = p(\mathbf{z}_{\text{sup}}|y)$ is the prior distribution of the superclass that y belongs to, and $p(x|\mathbf{z}_{\text{sup}}^y) = p(x|\mathbf{z}_{\text{sup}}, y)$ is the data distribution conditioned on the superclass $\mathbf{z}_{\text{sup}}^y$.

- *Programs*

BPL [34], [35], [36], [37] uses a Bayesian process to model the generation process of character objects as a probabilistic program. This program will experience a bottom-up parsing analysis of primitives, sub-parts, parts, types, tokens and images. Furthermore, the intermediate types and tokens within the generation program are treated as the latent variable \mathbf{z}_{prog} . With the explicit probabilistic program for each character concept, BPL is able to access the compositionality and causality of character objects and can perform one-shot

classification, generate new exemplars given one sample, and generate new character classes as well.

- *Splits*

Chopping model [61] introduces the random data splits of the auxiliary set D_A as the latent variable \mathbf{z}_{spl} to bridge the mathematical dependence between raw image x and the label y . It makes many splits on D_A by assigning label 1 to half of auxiliary classes and 0 to the others, and then trains a predictor for each split. For one image in D_T , Chopping model will combine the predictions from all split-specific predictors to achieve the Bayesian posterior decision.

- *Reconstruction*

Unlike BPL [34], [35], [36], [37], a compositional patch model (CPM) that does not rely on the knowledge of dynamic strokes in character images is proposed in [63]. Similar to BPL, the core assumption of this model is that the congenital character images share the same patch-based structure. Thus, this model first segments the single sample in D_{trn} for each class into a set of components, and then utilizes an AND-OR graph to reconstruct the test sample in D_{tst} . The reconstruction is essentially the latent variable \mathbf{z}_{rec} , which is used to make the final one-shot classification for test samples.

- *Statistics*

Neural Statistician model [65] deploys a deep network to produce statistics that encapsulate a generative model for each D_{trn} . Concretely, the statistics are described by a mean and variance specifying a Gaussian distribution in the latent space. Using the latent variable \mathbf{z}_{stat} , Neural Statistician can realize one-shot generation and classification.

Table 1 presents an intuitive comparison between the above mainstream generative model based FSL approaches that constructed the latent variable \mathbf{z} from different perspectives. Except for the Neural Statistician, the remainders were born in the non-deep period of FSL development process, and most of them are tailored in light of the specific task form or data form, lacking the scalability to more general cases. Besides, these early works were validated on various experimental datasets with different evaluation settings, having not formed some comparable benchmarks for subsequent FSL researches at that time.

4 DISCRIMINATIVE MODEL BASED APPROACHES

Unlike the above FSL approaches based on generative model, the discriminative model based FSL approaches attempt to model the posterior probability $p(y|x)$ directly for task T using the scarce training set D_{trn} . The computation model for $p(y|x)$ generally contains a feature extractor and a predictor. For few-shot image recognition tasks, for example, the feature extractor and the predictor respectively might be a CNN and softmax layer. Due to the sample scarcity in D_{trn} , it would be easy to trap into overfitting when fitting $p(y|x)$ only with D_{trn} . Therefore, existing discriminative model based FSL approaches pursue the construction of $p(y|x)$ from different perspectives. We summarize them into the following classes.

The first one is based on augmentation, which advocates learning a general augmentation function $\mathcal{A}(\cdot)$ from the auxiliary dataset D_A to augment the samples or the features of samples in D_{trn} . The augmentation based FSL approaches are reviewed in Section 4.1. The second one is based on metric learning, which aims to learn a pairwise similarity metric $S(\cdot, \cdot)$ on D_A . By this metric, a nearest-neighbor classifier can be used for final prediction. The metric learning based FSL approaches are introduced in Section 4.2. The third is based on meta learning, which leverages D_A to construct many tasks similar to the task T and adopts the cross-task training strategy to distill some transferrable models, algorithms or parameters. The meta learning based FSL approaches are detailed in Section 4.3. Besides, there also exist some other FSL approaches, which are discussed in Section 4.4.

4.1 Augmentation

Augmentation is an intuitive way to increase the number of training samples and enhance data diversity. In the field of vision, some basic augmentation operations include rotating, flipping, cropping, translation, and adding noise into images [11], [67], [68]. For FSL tasks, these low-level augmentation means are insufficient to bring essential gains in the generalization ability of FSL models. In this context, more sophisticated augmentation models, algorithms or networks customized for FSL were proposed, and they mainly occurred in the deep period. Fig. 2 illustrates the general framework of augmentation based FSL approaches. Except for DAGAN [69] that augments the samples in D_{trn} at the data level, other approaches achieve the feature-level augmentation for training samples in task T . According to whether their augmentation relies on external side information (such as semantic attributes [70], word vectors [71], etc), we further divide the existing augmentation based FSL approaches into supervised and unsupervised ones.

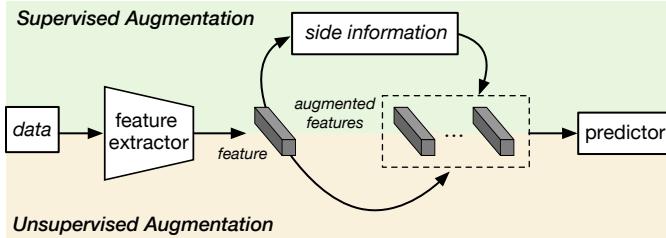


Fig. 2. General framework of augmentation based FSL approaches.

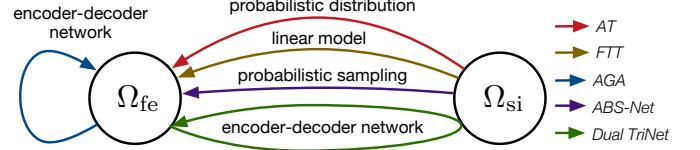


Fig. 3. Mapping relationship between the feature and side information in terms of different supervised augmentation approaches.

4.1.1 Supervised Augmentation

Several FSL approaches based on supervised augmentation include Feature Trajectory Transfer (FTT) [72], AGA [73], Dual TriNet [74], [75], Author-Topic (AT) [42] and ABS-Net [76]. For ease of notation, let Ω_{fe} be the feature space, and Ω_{si} be the side information space. The augmentation $\mathcal{A}(\cdot)$ learned by these approaches, essentially, is a mapping relationship between Ω_{fe} and Ω_{si} , although they differ in the mapping direction and mapping module, as shown in Fig. 3.

FTT [72] focuses on one-shot scene image classification, which leverages the consecutive attributes in scene images (e.g., “rainy”, “dark” or “sunny”) to directionally synthesize the features for the one-sample task scene class. In particular, FTT suggests to learn a linear mapping trajectory on auxiliary scene classes that maps attribute $a \in \mathbb{R}_+$ to feature $x \in \mathbb{R}^d$:

$$x = w \cdot a + b + \epsilon, \quad (7)$$

where $w, b \in \mathbb{R}^d$ are learnable parameters and ϵ denotes Gaussian noise. This mapping trajectory is expected to be transferrable from auxiliary classes to task classes. As shown in Fig. 4, given only one training sample for a task scene class, one can artificially set the strength of its attribute (e.g., the degree of *sunny*) to produce many synthetic features by the well-learned linear mapping trajectory in Eq. (7). However, we must note that FTT requires the fine-grained and consecutive attribute annotation, which is a prohibitive cost for data preparation.

Comparably, AGA [73] develops an encoder-decoder network to map the feature of a sample into another synthetic feature at a different attribute strength with the input feature. For example, as shown in Fig. 5, AGA aims to learn a class-agnostic feature transfer module $\phi_{[1,2]}^3$ on auxiliary classes (e.g., Tables, Chairs) that takes the features of objects with

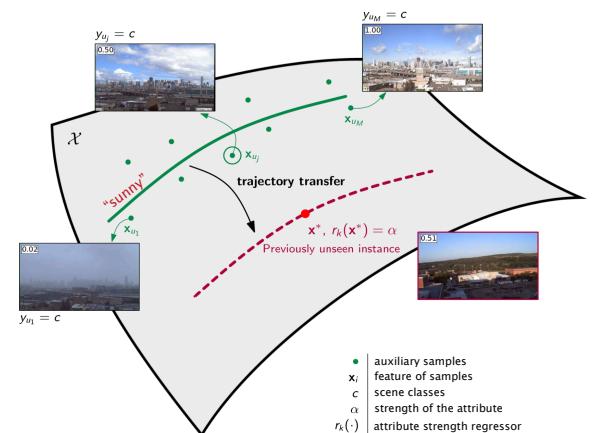


Fig. 4. Illustration of Feature Trajectory Transfer (FTT) [72]. Number at upper right corner in each scene image describes the strength of “sunny”.

TABLE 2
Summary of supervised (top part) or unsupervised (bottom part) augmentation based FSL approaches

Approaches	Side Information	Mapping Direction	Mapping Module	Task Type	Experimental Dataset
FIT [72]	transient attributes (<i>rainy, sunny, etc</i>)	$\Omega_{\text{si}} \rightarrow \Omega_{\text{fe}}$	linear model	scene location classification	Transient Attributes Database (TADB) [77], SUN Attributes Database (SADB) [78]
AGA [73]	attribute strength (<i>depth, pose</i>)	$\Omega_{\text{fe}} \xrightarrow{\Omega_{\text{si}}} \Omega_{\text{fe}}$	encoder-decoder network (MLP)	2D/3D object classification	SUN RGB-D [79]
AT [42]	discrete attributes (<i>black, fierce, etc</i>)	$\Omega_{\text{si}} \rightarrow \Omega_{\text{fe}}$	probabilistic distribution	image classification	Animals with Attributes (AwA) [80]
Dual TriNet [74], [75]	word vectors, discrete attributes	$\Omega_{\text{fe}} \rightarrow \Omega_{\text{si}} \rightarrow \Omega_{\text{fe}}$	encoder-decoder network (CNN)	image classification	miniImageNet [49], Cifar-100 [81], CUB [82], Caltech-256 [83]
ABS-Net [76]	discrete attributes (<i>ForColor, BackColor</i>)	$\Omega_{\text{si}} \rightarrow \Omega_{\text{fe}}$	probabilistic sampling	image classification	Colored MNIST [76]
GentleBoostKO [39]	-	$\Omega_{\text{fe}} \rightarrow \Omega_{\text{fe}}$	knockout (feature element replacement)	binary image classification	Caltech datasets [84]
SH [85]	-	$\Omega_{\text{fe}} \rightarrow \Omega_{\text{fe}}$	quadruplet-based MLP (3 features \rightarrow 1 feature)	image classification	ImageNet1k [23]
Hallucinator [86]	-	$\Omega_{\text{fe}} \rightarrow \Omega_{\text{fe}}$	MLP-based generator (1 features \rightarrow 1 feature)	image classification	ImageNet1k [23]
CP-ANN [87]	-	latent space $\rightarrow \Omega_{\text{fe}}$	GAN	image classification	ImageNet1k [23]
Δ -encoder [88]	-	$\Omega_{\text{fe}} \rightarrow \Omega_{\text{fe}}$	encoder-decoder network (MLP) (3 features \rightarrow 1 feature)	image classification	miniImageNet [49], Cifar-100 [81], CUB [82], Caltech-256 [83], AwA [80], aPascal&aYahoo (APY) [89]
DAGAN [69]	-	$\Omega_{\text{da}} \rightarrow \Omega_{\text{da}}$	GAN	image generation, image classification	Omniglot [37], EMNIST [90], VGG-Faces [91]
IDeMe-Net [92]	-	$\Omega_{\text{da}} \rightarrow \Omega_{\text{da}}$	Deformation Sub-network (2 images \rightarrow 1 images)	image classification	ImageNet1k [23], miniImageNet [49]

Note: the term in the column of "Task Type" without "binary" all indicates multi-class classification.

depth in the range of 1-2 [m] as inputs and outputs their synthetic features with the depth of 3 [m]. Using this feature transfer module, one can augment the single sample of one task class with various depth strength. This idea seems similar to FIT [72], but two different points exist between them. First, the mapping direction of AGA is $\Omega_{\text{fe}} \rightarrow \Omega_{\text{fe}}$, while that of FIT is $\Omega_{\text{si}} \rightarrow \Omega_{\text{fe}}$. Second, the feature synthesis in FIT is guided by directly allocating the desired attribute strength into its linear mapping model, but that in AGA is achieved by the encoder-decoder network specializing in the feature mapping between two explicit attribute strength.

Similarly, Dual TriNet [74], [75] also utilizes an encoder-decoder network to achieve the feature-level augmentation.

Apart from the difference in the architecture of encoder-decoder network used by them (Dual TriNet uses CNN, while AGA uses MLP), another major difference is the bottleneck embedding between the encoder and decoder: the bottleneck embedding of Dual TriNet is a semantic attribute or word vector, while that of AGA is trivial latent embedding. Dual TriNet models the bottleneck embedding space as a Semantic Gaussian or a Semantic Neighbourhood, in which large amounts of semantic vectors can be sampled to be decoded into synthetic features. From this point of view, the mapping direction of Dual TriNet is $\Omega_{\text{fe}} \rightarrow \Omega_{\text{si}} \rightarrow \Omega_{\text{fe}}$.

AT [42] uses a topic model [93] to model the relationship between images and attributes, where each image is treated as a document containing a mixture of topics (i.e., attributes), and each topic is represented by a probabilistic distribution of words (i.e., features). The parameters of this probabilistic distribution are estimated on the auxiliary dataset D_A . By the explicit distribution, a large amount of features of a specific class can be generated given the attributes of this class.

ABS-Net [76] first conducts an attribute learning process on the auxiliary dataset D_A , which allows the establishment of a repository of attribute features. Given the attribution description of one class, a probabilistic sampling operation is performed on the repository, which maps the attributes to the pseudo features of this class.

The top part of Table 2 summarizes the main characteristics of these supervised augmentation based FSL approaches. Considering the labeling cost of side information, the above approaches are more suitable for the task or dataset containing some side information.

4.1.2 Unsupervised Augmentation

Typical unsupervised augmentation based FSL approaches include GentleBoostKO [39], Shrinking and Hallucinating

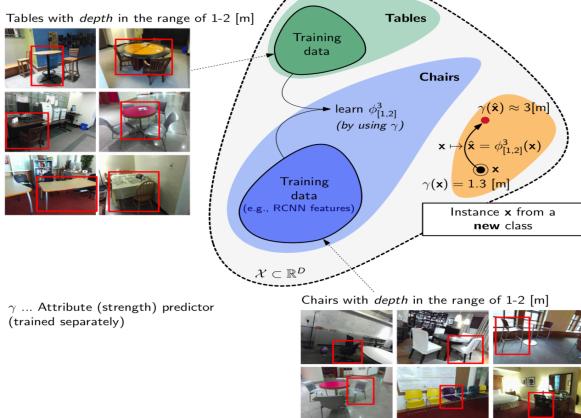


Fig. 5. Illustration of Attribute-Guided Augmentation (AGA) [73]. *Depth* is the attribute. Tables and Chairs are two auxiliary classes. $\phi_{[1,2]}^3$ is an encoder-decoder network that transfers the feature of an object with depth in the range of 1-2 [m] into another feature with depth of 3 [m].

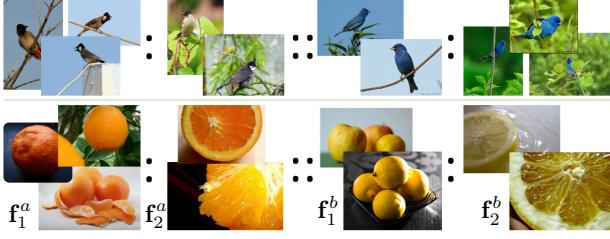


Fig. 6. Illustration of transformation analogies in the form of quadruplets used by SH [85]. Top row: birds with sky background *versus* birds with greenery background. Bottom row: whole fruits *versus* cut fruit.

(SH) [85], Hallucinator [86], CP-ANN [87], Δ -encoder [88], DAGAN [69] and IDeMe-Net [92] etc, which seek to augment the data or features without any external side information.

GentleBoostKO [39] is a straightforward FSL solution in the early non-deep period that synthesizes features by a knockout procedure. The knockout is realized by replacing one element of a feature with an element of another feature in the same coordinate. Its key insight is to create corrupted copies of the very few samples to increase the robustness.

SH [85] was built on the motivation that the intra-class variation can generalize across classes (e.g., pose transformations), which is similar to the intuition of FTT [72] and AGA [73]. The difference among them is that the intra-class variation in FTT [72] and AGA [73] can be explicitly described by the side information (e.g., the strength of *sunny* attribute, the *depth* value of object), while the underlying intra-class variation in SH needs to be mined from implicit transformation analogies in the form of quadruplets ($f_1^a, f_2^a, f_1^b, f_2^b$), as shown in Fig. 6, where a, b denote two classes. These quadruplets are mined from the auxiliary set by an unsupervised clustering and many heuristic steps. Furthermore, an MLP-based mapping module G can be learned based on these quadruplets, which takes three features as inputs and outputs a synthetic feature, i.e., $\hat{f}_2^b = G(f_1^a, f_2^a, f_1^b)$. Given only one training sample for a task class, one can deduce analogically other synthetic features for this class by G .

The high-level motivation of Δ -encoder [88] is similar to that of AGA [73], FTT [72] and SH [85]. It also suggests to extract transferrable intra-class variation (called Δ) from auxiliary set D_A and apply this variation to the novel task classes so as to synthesize new samples for the task classes. Similar to SH [85], Δ -encoder also transfers Δ based on the underlying quadruplet analogy, and the main difference is the specific mapping module that deals with the quadruplet

relationship: SH [85] uses a trivial MLP but Δ -encoder develops an encoder-decoder network whose bottleneck embedding is expected to capture the intra-class variation Δ .

As shown in Fig. 7, Hallucinator [86] uses an MLP-based generator G to augment features for the training samples in D_{trn} , i.e., $\hat{f} = G(f, z)$, where f is an original feature and z is a noise vector. This generator was devised into a plug-and-play module that can be incorporated into a variety of ready-made meta learning modules, such as Matching Nets [49], Prototypical Nets [54] or Prototype Matching Nets [86]. The meta learning FSL approaches will be reviewed in Section 4.3.

CP-ANN [87] achieved feature augmentation for the few support samples via a Generative Adversarial Networks (GAN) [94] based set-to-set translation, which was designed to preserve the covariance of auxiliary samples during augmentation. DAGAN [69] takes the samples in D_{trn} as input and generates the within-class data directly (i.e., $\Omega_{\text{da}} \rightarrow \Omega_{\text{da}}$) by a conditional GAN. Z. Chen *et al.* [92] insisted the visual fusion between two similar images may maintain critical semantic information and contribute to formulating the decision boundaries of the final classifier, and thus they proposed IDeMe-Net to generate the deformed images for the small amounts of support samples. Similar to Hallucinator [86], both DAGAN [69] and IDeMe-Net [92] were designed to work in coordination with other ready-made meta learning based FSL approaches like Matching Nets [49] and Prototypical Nets [54]. A high-level summary for the above unsupervised augmentation based FSL approaches is made in the bottom part of Table 2.

4.1.3 Discussion

We must emphasize that augmentation based FSL approaches do not conflict with other FSL approaches, such as those based on metric learning or meta learning to be discussed in Section 4.2 and 4.3. On the contrary, most of the augmentation based FSL approaches are complementary to them and they can be used as the plug-and-play module: one can first adopt these augmentation strategies to enrich D_{trn} and then learn on the augmented $D_{\text{trn}}^{\text{aug}}$ through other FSL approaches.

4.2 Metric Learning

The general objective of metric learning [45] is to learn a pairwise similarity metric $S(\cdot, \cdot)$ under which a similar sample pair can obtain a high similarity score while the dissimilar pair gets a low similarity score. All FSL approaches based on metric learning adhere to this principle, as shown in Fig. 8, which create the similarity metric using auxiliary dataset D_A and generalize it to the novel classes of task T . The similarity metric could be a simple distance measurement, a sophisticated network or other feasible modules or algorithms as long as they can estimate the pairwise similarity between

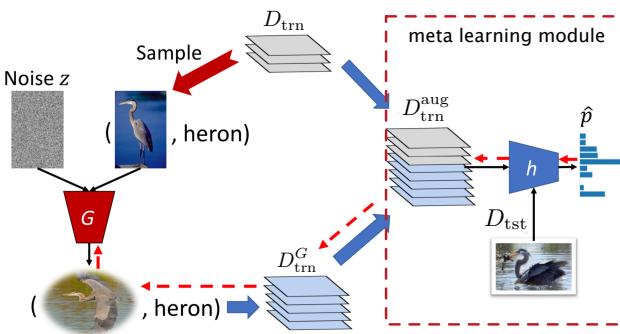


Fig. 7. Framework of Hallucinator [86]. D_{trn}^G : the augmented sample set.

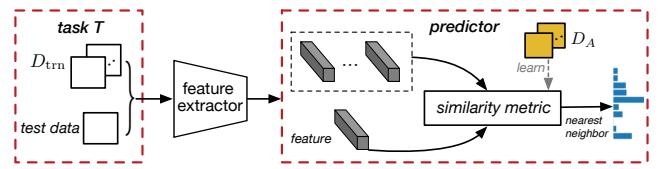


Fig. 8. General framework of metric learning based FSL approaches.

TABLE 3
Summary of metric learning based FSL approaches

Approaches	Similarity Metric $\mathcal{S}(\cdot, \cdot)$	Metric Loss	Task Type	Experimental Dataset
CRM [38]	$d(x_i, x_j)$ (Mahalanobis distance)	hinge loss	image classification	Latin Character database [38]
KernelBoost [95]	$K(x_i, x_j)$ (kernel function)	exponential loss	image classification, image retrieval	UIC [96], MNIST [60], YaleB [97]
Siamese Nets [30]	$\mathbf{p}(x_i, x_j)$ (siamese CNN)	binary cross-entropy loss	image classification	Omniglot [37]
Triplet Ranking Nets [98]	$d(x_i, x_j)$ (Euclidean distance)	triple ranking loss	image classification	Omniglot [37], miniImageNet [49]
SRPN [99]	$\mathbf{p}(x_i, x_j)$ (GAN+siamese CNN)	adversarial loss	image classification	Omniglot [37], miniImageNet [49]
MM [100]	$d(x_i, x_j)$ (memory+dot product)	memory loss	image classification, translation	Omniglot [37], WMT14 [100]
AdaptHistLoss [101]	$d(x_i, x_j)$ (cosine distance)	histogram loss	image classification, translation	MNIST [60], Isolet of UIC [96], Omniglot [37], tinyImageNet [101]

samples or features. Several representative metric learning based FSL approaches include Class Relevance Metrics (CRM) [38], KernelBoost [95], Siamese Nets [30], Triplet Ranking Nets [98], Skip Residual Pairwise Net (SRPN) [99], Memory Module (MM) [100] and AdaptHistLoss [101]. They developed various forms of similarity metrics associated with different metric loss functions to address FSL tasks.

CRM [38] is a foundation work of metric learning based FSL approach proposed in the non-deep period. It uses the Mahalanobis distance to measure the pairwise similarity:

$$d(x_i, x_j) = \sqrt{(x_i - x_j)^\top A(x_i - x_j)} = \|Wx_i - Wx_j\|_2, \quad (8)$$

where $A = W^\top W$ is a symmetric positive semi-definite matrix that needs to be learned from auxiliary dataset D_A . The learning objective of CRM follows the form of hinge loss to make the distance of positive sample pair (x_i, x_i^+) smaller than that of negative sample pair (x_j, x_j^-) by γ at least:

$$d(x_i, x_i^+) \leq d(x_j, x_j^-) - \gamma. \quad (9)$$

Once trained on D_A , the Mahalanobis distance is applied to the task T to enable the nearest neighbor (NN) classification.

KernelBoost [95] suggests to learn the pairwise distance in the form of a kernel function through a boosting algorithm. The kernel function is defined as a combination of some weak kernel functions, $K(x_i, x_j) = \sum_{t=1}^T \alpha_t K_t(x_i, x_j)$. Each weak kernel $K_t(\cdot, \cdot)$ learns a Gaussian Mixture Model (GMM) of the data, and $K_t(x_i, x_j)$ represents the probability that both x_i and x_j belong to a same Gaussian component within the t -th GMM. The kernel K is optimized by an exponential loss:

$$\ell = \sum_{i,j} \exp(-y_{ij}K(x_i, x_j)), \quad (10)$$

where y_{ij} is 1 if x_i and x_j are from the same class, and -1 otherwise. Finally, a kernel NN classifier can be formed.

Siamese Nets [30] is the first work that brings deep neural networks into FSL tasks. It consists of twin CNNs that share the same weights. The twin CNNs accept a pair of samples

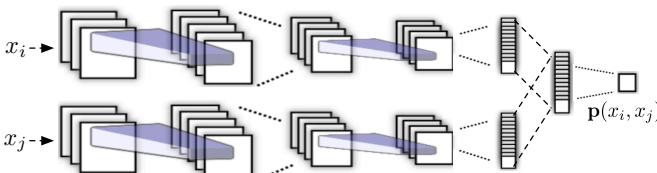


Fig. 9. Architecture of Siamese Nets [30]. Twin CNNs share weights.

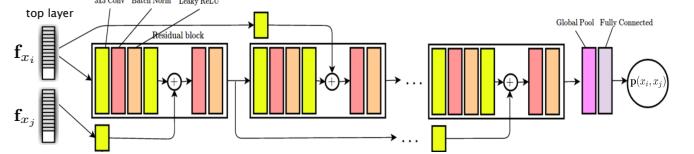


Fig. 10. Architecture of Skip Residual Pairwise Net (SRPN) [99].

(x_i, x_j) as inputs and their outputs at the top layer are combined in order to output a single pairwise similarity score $\mathbf{p}(x_i, x_j)$, as depicted in Fig. 9. The twin CNNs are trained through the following binary cross-entropy loss:

$$\ell = \sum_{i,j} y_{ij} \log \mathbf{p}(x_i, x_j) + (1-y_{ij}) \log(1-\mathbf{p}(x_i, x_j)), \quad (11)$$

where $y_{ij}=1$ when x_i and x_j belong to the same class, and 0 otherwise. The well-trained twin CNNs are frozen and used as a fixed similarity metric to make inference on FSL task T in the manner of NN.

Triplet Ranking Nets [98] extended Siamese Nets [30] from pairwise samples to triplets and used the triplet ranking loss [102] to optimize the metric space and then capture the similarity between samples. SRPN [99] is also an evolution of Siamese Nets [30], which involves two main modifications: (1) Replacing the simple pairwise combination at top layer used by Siamese Nets with a more sophisticated skip residual network [13] that separates the intermediate computations for the pair of samples, as shown in Fig. 10. (2) Using an additional GAN [94] to regularize the skip residual network by taking the skip residual network as a discriminator network and introducing another auto-encoder based generator. Thus, the metric loss of SRPN is naturally incorporated into the adversarial loss of GAN.

MM [100] develops a life-long memory module to learn the similarity metric, which regards the feature of test data as its query q and stores many continuously updated keys associated with values (*i.e.*, class labels). This memory is optimized by the following hinge-based memory loss:

$$\ell = \sum_q [q^\top k^- - q^\top k^+ + \gamma]_+, \quad (12)$$

where k^- and k^+ are respectively the negative and positive key in terms of q . Whether a key is positive or negative is determined by comparing its value to the class label of q .

AdaptHistLoss [101] advise to adopt histogram loss [103] to learn a feature space where the simple cosine distance can effectively measure the similarity between two features. Histogram loss suggests to construct two sets of similarities,

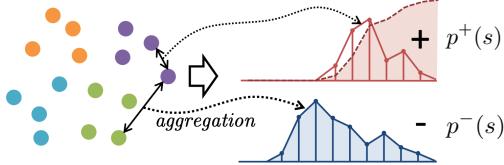


Fig. 11. Computation of histogram loss [103] for a batch of samples. Dots denote the features of samples. The same color indicates the same class.

$S^+ = \{s(\mathbf{f}_{x_i}, \mathbf{f}_{x_j}) | y_i = y_j\}$ and $S^- = \{s(\mathbf{f}_{x_i}, \mathbf{f}_{x_j}) | y_i \neq y_j\}$, where \mathbf{f}_{x_i} denotes the feature of data x_i with class label y_i and $s(\cdot, \cdot)$ is the feature-level similarity metric (i.e., cosine similarity). Using S^+ and S^- , one can estimate the similarity distributions of positive and negative pairs as the histograms, as shown in Fig. 11, which are denoted as $p^+(s)$ and $p^-(s)$ respectively. Then, the histogram loss is defined as the reverse probability that the similarity in a random negative pair is more than the similarity in a random positive pair:

$$\ell = \int_{-1}^1 p^-(s) \left[\int_{-1}^s p^+(z) dz \right] ds = \mathbb{E}_{s \sim p^-} \left[\int_{-1}^s p^+(z) dz \right]. \quad (13)$$

Since histogram loss only focuses on the similarity distributions of positive pairs and negative pairs but agnostic to class labels, this metric can be directly transferred to FSL task T .

Table 3 summarizes the main characteristics of the above metric learning based FSL approaches. In addition to these approaches, it should be noted that the Learn-to-Measure FSL approaches (see Section 4.3.1), strictly speaking, all belong to the scope of metric learning. Considering they pursue the similarity metric under the paradigm of meta learning, we discuss them in the section of meta learning.

4.3 Meta Learning

The idea of meta learning was proposed as early as the 1990s [104], [105], [106]. As deep learning grew in popularity, some works proposed to utilize the meta learning policy to learn to optimize deep models [107], [108], [109]. In general, meta learning advocates to learn across tasks and then adapt to new tasks, as shown in Fig. 12, which aims to learn on the level of tasks instead of samples, and learns the task-agnostic learning systems instead of task-specific models.

FSL is a natural testbed to validate the capability of meta learning approaches across tasks where only a few labeled samples are given per task. Meta learning approaches process FSL problems in two stages: meta-train and meta-test. In meta-train, the model is exposed to many independent supervised tasks $T \sim p(T)$ that are constructed on the auxiliary dataset D_A (also called “episode” [49], [50]) to learn how to adapt to future related tasks, where $P(T)$ defines a task distribution and the word of *related* means that all tasks

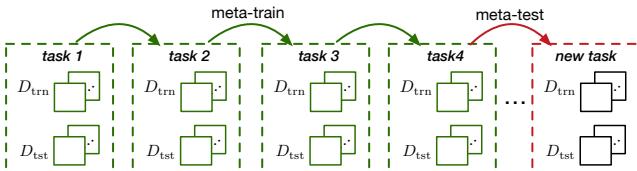


Fig. 12. General framework of meta learning based FSL approaches.

are from $P(T)$ and follow the same task paradigm, e.g., all tasks are C -way K -shot problems. Each meta-train task T entails a task-specific dataset $D_T = \{D_{trn}, D_{tst}\}$, where $D_{trn} = \{(x_i, y_i)\}_{i=1}^{N_{trn}}$ and $D_{tst} = \{(x_i, y_i)\}_{i=1}^{N_{tst}}$. In meta-test, the model is tested on a new task $T \sim p(T)$ whose label space is disjoint with the labels seen during meta-train. In most cases, D_{trn} is called as support or description set and D_{tst} is called as query set. Accordingly, the samples of them are called as support samples and query samples, respectively. The meta learning objective is to find the model parameters θ that minimize the expected loss $L(\cdot; \theta)$ across all tasks:

$$\min_{\theta} \mathbb{E}_{T \sim P(T)} L(D_T; \theta). \quad (14)$$

We must emphasize that the meta learning is a high-level cross-task learning strategy rather than a specific FSL model. Based on what the meta learning model seeks to meta-learn behind this learning strategy, we generally summarize the meta learning based FSL methods into five sub-categories: Learn-to-Measure (L2M), Learn-to-Finetune (L2F), Learn-to-Parameterize (L2P), Learn-to-Adjust (L2A) and Learn-to-Remember (L2R).

4.3.1 Learn-to-Measure

The L2M approaches inherit the main idea of metric learning in essence as shown in Fig. 8, but they are different from the metric learning based FSL approaches as described in Section 4.2 in the implementation level: the L2M approaches adopt the meta learning policy to learn the similarity metric that is expected to be transferrable across different tasks. L2M has always been an important branch of meta learning based FSL approaches, and several milestone meta learning approaches such as Matching Nets [49], Prototypical Nets [54], and Relation Net [55] all belong to the L2M category.

We first describe the general pipeline of L2M mathematically. For a task T , let x_i be a support sample in D_{trn} and x_j be a query sample in D_{tst} , and let $f(\cdot; \theta_f)$ and $g(\cdot; \theta_g)$ be the embedding models that map the support and query samples into features respectively. Moreover, all L2M approaches contain a metric module $S(f, g; \theta_S)$ to measure the similarity between support and query samples, which might be a parameter-free distance metric (e.g., Euclidean distance, cosine distance) or a learnable network. The similarity output by this metric module is used to form the final predicted probability of the query sample. Existing L2M approaches are different mainly in the model design and selection of f , g and S , and we draw the development relationship between different L2M approaches in Fig. 13 in a highly abstract way.

- *Prototypical Nets and its Variants*

The pioneer of L2M is Micro-set Learning [43], although the concept of meta learning was not mentioned by it at that time. This approach artificially constructs many micro-sets like test scenarios from the auxiliary dataset D_A , and each micro-set contains several support and query samples belonging to a few non-task classes. Both the embedding models f and g are realized through a weight-shared linear projection (i.e., $f = g$) and the similarity metric S is achieved by Euclidean distance. Moreover, NCA [110] is used to measure the final probability. Actually, the micro-sets are equivalent to the so-called *episodes* nowadays, and each micro-set is a meta-train task T . Importantly, if we replace the linear projection

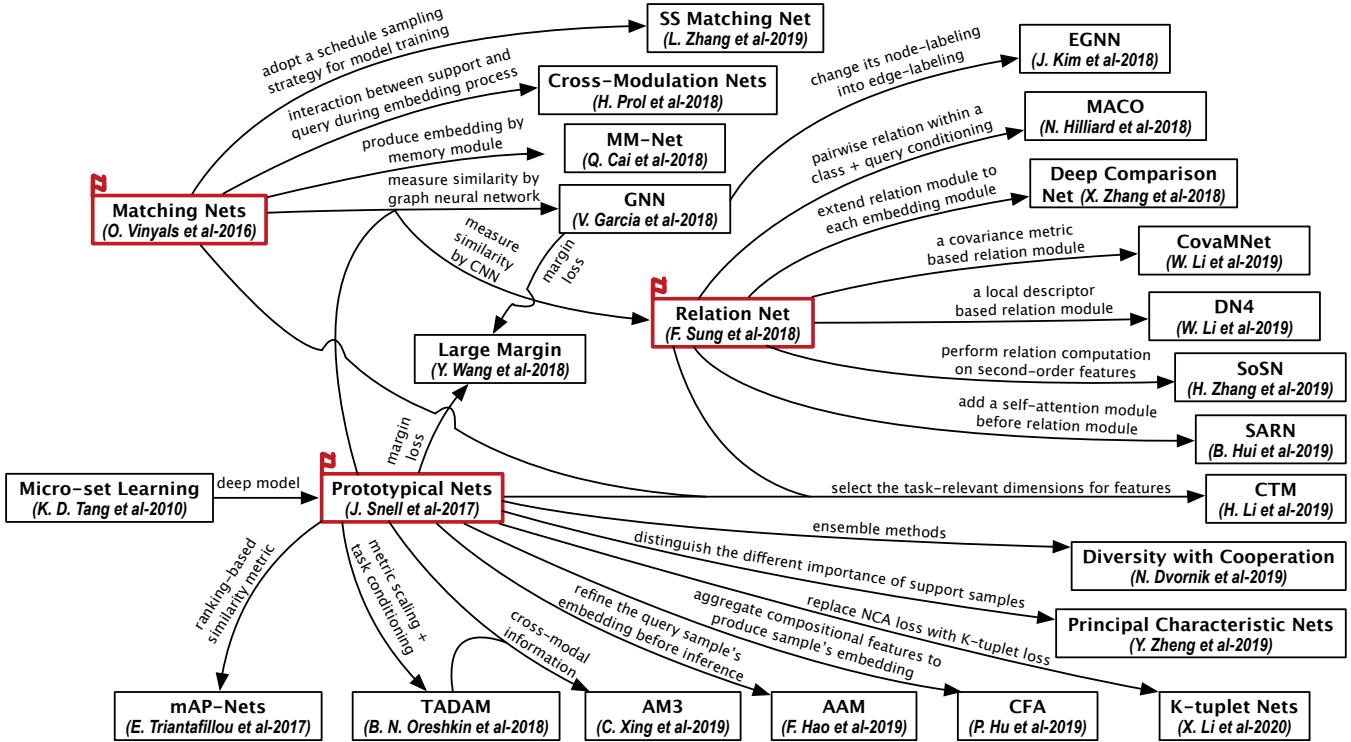


Fig. 13. Development relationship between different Learn-to-Measure FSL approaches.

model with a deep learning based embedding model such as CNN, Micro-set Learning [43] would evolve into the classic Prototypical Nets [54]. It takes the center of congener support samples' embeddings as the prototype of this class

$$p_c = \frac{1}{K} \sum_{(x_i, y_i) \in D_{\text{trn}}} \mathbb{1}(y_i == c) f(x_i; \theta_f), \quad (15)$$

and then also leverages the Euclidean distance based NCA like Micro-set Learning [43] to predict the probability:

$$P(y_j = c|x_j) = \frac{\exp(-d(g(x_j; \theta_g), p_c))}{\sum_{c'=1}^C \exp(-d(g(x_j; \theta_g), p_{c'}))}, \quad (16)$$

where f and g are also weight-shared embedding models (i.e., $f = g$). This L2M framework is an important cornerstone of many subsequent FSL approaches. In [111], mAP-Nets were proposed to learn an informative similarity metric from the perspective of information retrieval. It chooses to optimize an mAP-based ranking loss within each meta-train task using Structure SVM [112] or Direct Loss Minimization [113]. TADAM [114] further optimized the similarity metric S of Prototypical Nets by introducing a metric scaling factor α and transformed the original task-irrelevant f into a task-conditioning embedding model through a task embedding network (TEN) [114]. The TEN follows the key idea of FILM conditioning layer [115], and customizes some adjustment parameters (e.g., scaling and shift meta parameters) for the embedding model f in light of the current task representation. AM3 [116] incorporate extra cross-modal information (e.g., semantic representations) into Prototypical Nets and TADAM to enhance the metric learning process. Specifically, it used GloVe [117] to extract word embeddings for the semantic class labels and then built a new prototype by a convex

combination of both the visual feature and word embedding. AAM [118] proposed to refine the query sample's embedding before Eq. (16) to render it closer to its corresponding class center. CFA [119] achieved a compositional feature extraction for images instead of the vanilla image-to-vector mapping. K-tuplet Nets [120] changed the NCA loss of Prototypical Nets into a K-tuplet metric loss. Y. Zheng *et al.* [121] believed that the average prototype ignores the different importance of different support samples and thus proposed Principal Characteristic Nets. Diversity with Cooperation [122] achieved an ensemble of Prototypical Nets to encourage all individual networks to cooperate while encourage prediction diversity.

• Matching Nets and its Variants

The first deep learning based L2M approach is Matching Nets [49]. As shown in Fig. 14, it predicts the probability of query sample x_j by measuring the cosine similarity between the embedding of x_j and each support sample's embedding:

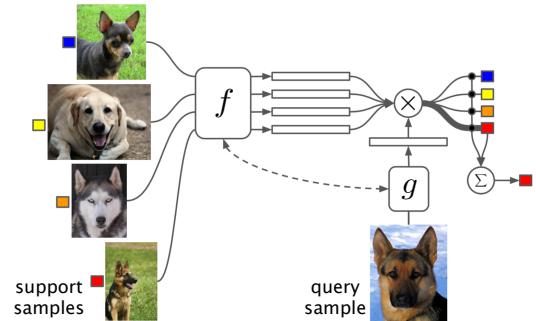


Fig. 14. Matching Nets architecture [49] (4-way 1-shot task for example).

$$p(\hat{y}_j|x_j, D_{\text{trn}}) = \sum_{(x_i, y_i) \in D_{\text{trn}}} a(x_j, x_i) \cdot \mathbf{y}_i, \quad (17)$$

where \mathbf{y}_i is a C -dimensional one-hot label vector (C -way FSL task) corresponding to the label $y_i \in \mathbb{R}^1$, and

$$a(x_j, x_i) = \frac{\exp(c(g(x_j; \theta_g), f(x_i; \theta_f)))}{\sum_{(x,y) \in D_{\text{trn}}} \exp(c(g(x_j; \theta_g), f(x; \theta_f)))}. \quad (18)$$

Matching Nets is different from Prototypical Nets in two aspects. First, the embedding model f and g of Matching Nets are two different networks. Concretely, f is a combination of CNN and BiLSTM [22] that aims to achieve full context embeddings (FCE) [49] for the few support samples, while g is a f -conditioning model that generates the features for query samples by a content attention mechanism. Second, the similarity metric S in Matching Nets is cosine distance instead of Euclidean distance. Several follow-up works have made some modifications and extensions based on Matching Nets. For example, Cross-Modulation Nets [123] modified the conditioning mechanism between f and g into a cross-modulation mechanism using FILM layers [115], which allows support and query samples to interact during the feature embedding process. MM-Net [124] developed a memory module [125] to produce feature embeddings, and the parameters of g are generated by this memory module. SS Matching Nets [126] considered the semantic diversity and similarity of class labels and exploited a scheduled sampling strategy to facilitate the model training of Matching Nets.

- *Relation Net and its Variants*

Unlike Prototypical Nets [54] and Matching Nets [49] which use the non-parametric Euclidean distance or cosine distance to measure the similarity between pairwise features, Relation Net [55] adopted a learnable CNN (denoted by $h(\cdot; \theta_h)$ here) to measure pairwise similarity, which takes the concatenation of feature maps of support sample x_i and query sample x_j as input and outputs their relation score $r(x_i, x_j)$, as shown in Fig. 15, which is formulated as

$$r(x_i, x_j) = h(C(f(x_i; \theta_f), g(x_j; \theta_g)); \theta_h) \in [0, 1], \quad (19)$$

where $f = g$ and C denotes the feature maps concatenation. It needs to be noted that, for Relation Net, the embeddings output by f (or g) are feature maps rather than feature vectors. Based on Relation Net, MACO [127] designed a relational stage after f to form the pairwise relation features within one

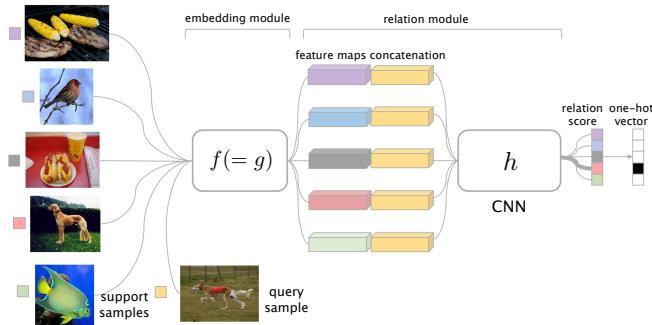


Fig. 15. Relation Net architecture [55] (5-way 1-shot task for example).

class and then used a query conditioning operation to predict the probability of query samples. Deep Comparison Net [128] extended Relation Net by deploying the relation module to each layer of the embedding model f . CovaMNet [129] and DN4 [130] replaced the relation module in Relation Net [55] by a covariance metric network and a deep local descriptor based image-to-class metric module, respectively. SoSN [131] chose to perform the relation computation on the second-order representation of feature maps. SARN [132] introduced a self-attention mechanism into Relation Net for capturing non-local features and enhancing representation.

- *Other L2M Approaches*

In addition to the above three mainstreams, there have also several hybrid variants based on them. For example, GNN [133] replaced the Euclidean distance in Prototypical Nets [54] and the cosine distance in Matching Net [49] with a learnable graph neural network where the nodes are set to the samples' embeddings and the edges are treated as the similarity between two samples. Conversely, EGNN [134] exchanged the roles of both nodes and edges in GNN [133], transforming it from a node-labeling framework to an edge-labeling one. Y. Wang *et al.* [135] proposed to impose a large margin constraint and augment the final classification loss of Prototypical Nets or GNN with a margin loss such that the metric space could be more discriminative. In consideration of the different importance of one feature element in different tasks, H. Li *et al.* introduced Category Traversal Module (CTM) [136] to select the most task-relevant dimensions of feature embeddings. CTM [136] can be used as a plug-and-play module for other approaches like Matching Nets [49], Prototypical Nets [54] and Relation Net [55].

4.3.2 Learn-to-Finetune

L2F approaches suggest to finetune a base learner for task T using its few support samples and make the base learner converge fast on these samples within several parameter update steps. Generally, every L2F approach contains a base learner and a meta learner. The base learner is for a specific task, which takes the sample as input and outputs the prediction probability. The base learner is learned by the higher-level meta-learner that is learned on a bunch of meta-train tasks to maximize the combined generalization power of the base learner on all tasks. Let θ_b and θ_m denote the parameters of base learner and meta learner, respectively. The learning process of L2F occurs at two levels. Gradual learning is performed across tasks, which aims to optimize the meta learning parameters θ_m and then facilitate the rapid learning of base learner for each specific task. Two milestone L2F approaches are MAML [50] and Meta-Learner LSTM [51].

MAML [50] is an elegant meta learning framework with strong interpretability, which has a profound influence on the field of meta learning and FSL. Its core idea is to search for a good parameter initialization for θ_b by cross-task training strategy such that the base learner with this initialization can rapidly generalize new tasks using a few support samples. Concretely, as the base learner copes with a task T , the one-step updated parameter θ'_b of base learner is computed as

$$\theta'_b = \theta_b - \alpha \nabla_{\theta_b} L(D_{\text{trn}}^T, \theta_b), \quad (20)$$

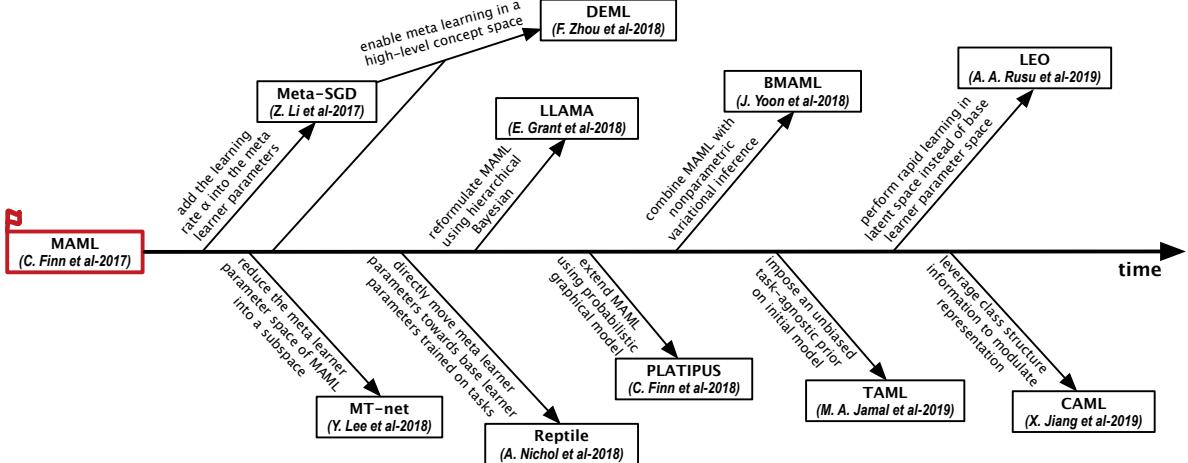


Fig. 16. Development relationship between MAML [50] and its variants.

where α is learning rate and $L(D_{\text{trn}}^T, \theta_b)$ is the loss on support set of task T when base learner parameter starts with θ_b . On the meta level, MAML optimizes meta learner by balancing the loss with updated base learner θ_b^T over many tasks:

$$\theta_m = \theta_m - \beta \nabla_{\theta_m} \sum_{T \sim P(T)} L(D_{\text{trn}}^T, \theta_b^T). \quad (21)$$

Note that the meta learner in MAML [50] is actually the base learner, that is, the meta learner parameter satisfy $\theta_m = \theta_b$. Eq. (20) is the rapid learning process that aims to finetune the base learner towards the specific task, while Eq. (23) is the gradual learning process that is intended to distill an appropriate parameter initialization for base learner.

Additionally, many L2F approaches belonging to MAML variants [137], [138], [139], [140], [141], [142], [143], [144], [145], [146] have been developed recently. The relationship between MAML and them is shown in Fig. 16. Meta-SGD [137] proposed to meta-learn not just the base learner initialization, but also the base learner update direction and learning rate. Thus, Meta-SGD modified the learning rate α in Eq. (20) into a learnable vector α and adds it into meta learner parameters:

$$\begin{aligned} \theta_b^T &= \theta_b - \alpha \circ \nabla_{\theta_b} L(D_{\text{trn}}^T, \theta_b), \\ (\theta_b, \alpha) &= (\theta_b, \alpha) - \beta \nabla_{(\theta_b, \alpha)} \sum_{T \sim P(T)} L(D_{\text{trn}}^T, \theta_b^T). \end{aligned} \quad (22)$$

Follow this line, DEML [140] made an incremental change for Meta-SGD, which equipped the meta-learner of Meta-SGD with a concept generator to enable learning to learn in a high-level concept space. In contrast, MT-net [138] proposed to reduce the meta learner parameter space of MAML into a subspace that is composed of each layer’s activation space and perform the rapid learning on this subspace. To avoid the computation of second-order derivative in MAML during gradual learning, A. Nichol *et al.* developed Reptile [139] that directly moves the meta learner parameter θ_m towards the base learner parameters θ_b^T that are updated on many tasks:

$$\theta_b = \theta_b - \beta \sum_{T \sim P(T)} (\theta_b - \theta_b^T). \quad (23)$$

LLAMA [141] reformulated MAML by hierarchical Bayesian and made an extension to MAML from the perspective of Bayesian posterior estimation. In consideration of the issue of task ambiguity when learning from small amounts of

samples, PLATIPUS [142] extended MAML using probabilistic graphical model and reframes it as a graph model inference problem, enabling simple and effective sampling of base learners for new tasks at meta-test time. In contrast, BMAML [143] coped with the model uncertainty when learning from a few samples by combining MAML with a non-parametric variational inference, i.e., Stein Variational Gradient Descent (SVGD) [147]. In addition, BMAML proposed a novel Chaser loss [143] during gradual learning to optimize the meta learner parameters θ_m . TAML [144] imposed an unbiased task-agnostic prior, which is achieved by an entropy-maximization/reduction or an inequality-minimization, on the initial model to prevent it from overperforming on meta-train tasks. LEO [145] designed a latent embedding \mathbf{z} , which is produced from support set of task T via an encoder $\mathbf{z} = E(D_{\text{trn}}^T; \theta_E)$, to generate the base learner parameter θ_b , i.e., $\theta_b = G(\mathbf{z}; \theta_G)$, and performed the rapid learning in the low-dimensional latent space instead of the high-dimensional base learner parameter space like Eq. (20):

$$\mathbf{z}^T = \mathbf{z} - \alpha \circ \nabla_{\mathbf{z}} L(D_{\text{trn}}^T, G(\mathbf{z}; \theta_G)). \quad (24)$$

Obviously, the combination of encoder $E(\cdot; \theta_E)$ and generator $G(\cdot; \theta_G)$ plays the role of meta learner of LEO, and thus its gradual learning process can be described as

$$\theta_m = \theta_m - \beta \nabla_{\theta_m} \sum_{T \sim P(T)} L(D_{\text{trn}}^T, G(\mathbf{z}^T; \theta_G)), \quad (25)$$

where the meta learner parameter $\theta_m = (\theta_E, \theta_G, \alpha)$. Compared with MAML, CAML [146] leverages the label structure to modulate the representations of base learner in light of the current task. Specifically, a parametric conditional transformation module is designed to perform the representation modulation. The combination of this module and the base learner acts as the meta learner of CAML, which is meta-learned via the gradual learning strategy adopted by MAML.

Another representative L2F approach is Meta-Learner LSTM [51], which suggests finetuning the base learner on the few support samples by a LSTM-based meta learner. As shown in Fig. 17, the LSTM-based meta-learner takes as input the loss and gradient of base learner with respect to each support sample, and its hidden state is treated as the updated base learner parameter, which would be used to handle the

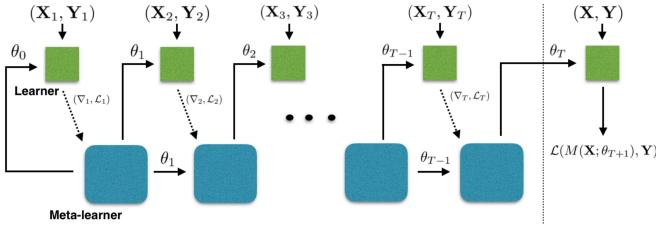


Fig. 17. Forward computational process of Meta-Learner LSTM for one task [51]. Green box is base learner and blue box is meta learner.

next support sample. In this framework, the vanilla gradient-based optimization for base learner parameters is superseded by an LSTM in the hope of learning appropriate parameter updates specifically for the scenario where a few updates will be made. J. Nie *et al.* extended Meta-Learner LSTM to a dual version, called Meta-Learner Dual-LSTM [148], and applied it to 3D model few-shot classification tasks. In addition, a recent L2F approach is MTL [149], which developed a light-weight scaling and shifting network attached to the frozen base learner to reduce the probability of overfitting when finetuned on the few support samples.

4.3.3 Learn-to-Parameterize

L2P is another kind of popular meta learning based FSL approach, which adheres to a straightforward idea: parameterizing the base learner or some subparts of base learner for a novel task so that it can address this task specifically. As shown in Fig. 18, most L2P approaches also contain both base learner and meta learner like L2F approaches, but the difference is that, for L2P approaches, the two learners are trained synchronously within each task and the meta learner is essentially a task-specific parameter generator. For a task T , the meta learner is expected to generate some T -specific parameters for the base learner or its subparts in light of the few support samples of task T and the state of current base learner that is handling these support samples. At this point, L2P approaches attempt to learn how to parameterize the base learner to render it applicable to the specific task.

Many L2P approaches have been proposed in recent years, which developed various parameter generation modules to parameterize different parts of base learner. They might parameterize the task-specific predictor in base learner [152], [153], [154], [157], [158], [159], [160], or the intermediate feature extraction layers in base learner [56], [150], [155], even or the whole base learner [151], [156]. The characteristics of different L2P approaches are summarized in Table 4.

Siamese Learnet [150], as shown in Fig. 19, used Siamese Nets [30] as its base learner where one intermediate convolutional layer (red block in Fig. 19) is designed to be dynamic to different tasks. Another single-stream Siamese

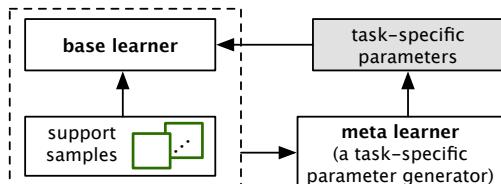


Fig. 18. General framework of Learn-to-Parameterize FSL approaches.

TABLE 4
Summary of Learn-to-Parameterize FSL approaches.

Approaches	Parameter Generation Module	Generated Parameters
Siamese Learnet [150]	single-stream Siamese Nets [30]	Conv. layer weights
Regression Nets [151]	MLP-based weight transformation	SVM weights
Dynamic Nets [152]	attention-based weight composition	predictor weights
Acts2Params [153]	MLP-based parameter predictor	Softmax layer weights
Imprinting [154]	MLP-based weight transformation	predictor weights
DCCN [155]	LSTM embedding module	Conv. layer weights
MeLA [156]	Auto-Encoder	all Conv. layer weights
DAE [157]	graph neural network	predictor weights
VERSA [158]	probabilistic amortization network	Softmax layer weights
R2-D2 [159]	ridge regression layer	predictor weights
MetaOptNet [160]	SVM	predictor weights
LGM-Net [56]	VAE-like weight generator	Conv. layer weights

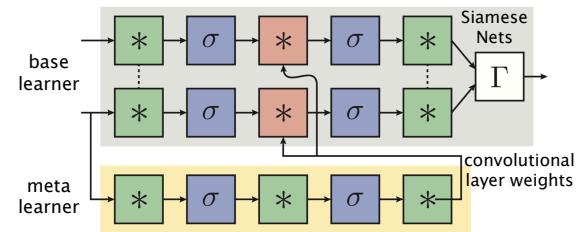


Fig. 19. Siamese Learnet architecture [150]. “*”: convolutional layer.

Nets is deployed as its meta learner to generate task-specific weights for this convolutional layer. LGM-Net [56] is a state-of-the-art L2P approach, which developed a MetaNet Module (i.e., meta learner) to generate the weights of TargetNet Module (i.e., base learner) on the basis of the few support samples in each FSL task, as shown in Fig. 20. Specifically, the MetaNet Module in LGM-Net takes the average embedding of support samples as input and produces the weights for each convolutional layer in base learner through an encoder-decoder model with multivariate Gaussian sampling. Once parameterized, the base learner of LGM-Net would make FSL inference similar to the classic Matching Nets [49].

Regression Nets [151] pursued a task-agnostic transformation of base learner's weights from a small-sample model to a large-sample model. Through this weight transformation, one can obtain more general weights for base learner albeit only on a few training samples. Dynamic Nets [152] advocated parameterizing the task-specific predictor by combining the average representation of the few support

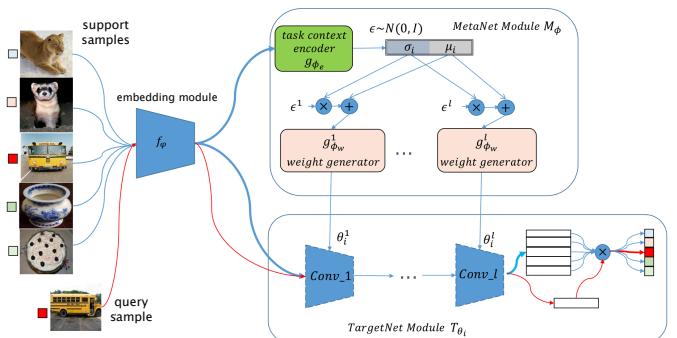


Fig. 20. LGM-Net architecture [56] (5-way 1-shot task for example).

samples and an attention-based weight composition on non-task predictor weights. Acts2Params [153] learned an MLP-based parameter predictor that maps neuron activations into the weights of the final Softmax predictor. Once well trained, it could directly predict the task-specific Softmax weights by taking the activations of the few support samples of this task as its inputs. Similarly, Imprinting [154] also inherited the mapping idea which transforms the embeddings of support samples into the task-specific predictor weights via an MLP. VERSA [158] exploited a versatile amortization network that accepts support samples as input and outputs the parameter distribution for task-specific Softmax predictor. DAE [157] used a graph neural network based denoising Auto-Encoder (AE) to generate final predictor parameters. Comparably, R2-D2 [159] adopted a differentiable ridge regression layer to parameterize the task-specific predictor, while MeteOptNet [160] advocated a differentiable convex optimization on SVM for generating final predictor weights. MeLA [156] is similar to LGM-Net [56] since they both tried to customize the convolutional layers of base learner (MeLA parameterized several rear convolutional layers while LGM-Net parameterized all convolutional layers) for the specific task via an encoder-decoder based generator (MeLA used AE while LGM-Net used an AE variant like Variational Auto-Encoder [161]).

4.3.4 Learn-to-Adjust

As depicted in Fig. 21, the core ideology of L2A approaches is to adaptively adjust the computation flow or computing nodes in the base learner for a specific sample to make this sample compatible with the base learner. One may find that both L2P and L2A are similar since they all use the meta learner to change the base learner, but L2A approaches have two distinctive characteristics different from L2P. (1) The degree of change on the base learner by L2A approaches is lighter since they only make some incremental adjustments to base learner instead of a complete parameterization to the base learner or its subparts like L2P. (2) The change on the base learner by L2A approaches is more fine-grained since the adjustment of L2A is sample-specific while the parameterization of L2P is task-specific.

Several typical L2A approaches include MetaNet [53], CSNs [162], MetaHebb [163] and FEAT [164], which differ in the selection for the parts needing to be adjusted as well as the design for the generated adjustment, as described in Table 5. MetaNet [53] deployed a fast-weight layer attached to each layer of the base learner. The weights of each fast-weight layer are meta-generated by an external meta learner in light of the input sample. These collateral branch layers are used to adjust the intermediate values of the input sample during the feedforward process. CSNs [162] selected to adjust the neuron state (i.e., pre-activation) of each hidden node in the base learner. Specifically, it combined a memory

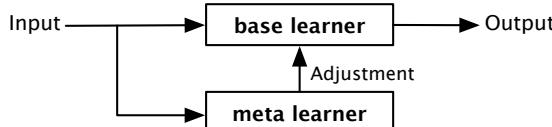


Fig. 21. General framework of Learn-to-Adjust FSL approaches.

TABLE 5
Summary of Learn-to-Adjust FSL approaches.

Approaches	Parts to be Adjusted	Adjustment
MetaNet [53]	parametric layers of base learner	layer-wise fast weights
CSNs [162]	neuron state of base learner	neuron-wise conditional shift
MetaHebb [163]	pre-Softmax layer of base learner	pre-Softmax fast-weight matrix
FEAT [164]	embeddings of support samples	vector-to-vector transformer

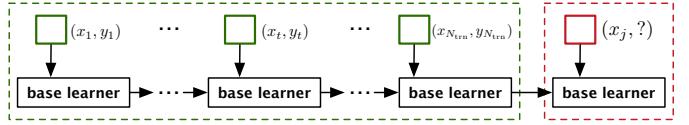


Fig. 22. General framework of Learn-to-Remember FSL approaches.

module with an attention-based memory read mechanism to generate the condition shift for each neuron pre-activation in the base learner. MetaHebb [163] added an auxiliary fast-weight matrix in pre-Softmax layer, which is meta-generated via Hebbian learning [165] and is expected to adjust the input's internal representation fed into final Softmax layer. FEAT [164] proposed to adapt the embeddings of support samples to render them more discriminative for the task in hand, and four kinds of set-to-set functions including BiLSTM [22], DeepSets [166], GCN [167] and Transformer [168] that aim to transform the original embedding vector to the adapted vector were investigated by [164].

4.3.5 Learn-to-Remember

Several representative FSL approaches, such as MANN [169], ARCs [170], SNAIL [171] and APL [172], belong to the kind of L2R. As shown in Fig. 22, its primary idea is to model the support set of an FSL task as a sequence and formulate the FSL task as a sequence learning task, where the query sample is required to match with previously seen information (i.e., support samples). Thus, the base learner of L2R approaches usually entails a temporal network to handle the few support samples. For example, MANN [52] utilized a memory-augmented Neural Turing Machine (NTM) [173] to rapidly assimilate the support samples and then retrieve them when the query sample arrives. ARCs [170] developed an attention based RNN to realize dynamic comparison between samples. SNAIL [171] devised a temporal convolutional network with soft attention to aggregate previously seen information and pinpoint specific information. APL [172] designed a surprised-based memory network to remember the most informative support samples it has encountered.

4.3.6 Discussion

The above five kinds of meta learning approaches all focus on dealing with FSL problems, however, each of them has their strengths or weaknesses. The L2M approaches would not be restricted by the specific settings of test scenarios since they only leverage the similarity between samples to make ultimate inference regardless of the number of classes and support samples per class (i.e., way/shot-agnostic). The L2F approaches need to be finetuned on each new task using the few support samples, which may yield a relatively long adaptation period to prepare for each task. One common

challenge faced by L2P and L2A approaches is a large number of model parameters because they have to deploy another meta learner completely different from base learner to generate a series of model parameters or adjustment parameters. Besides, the model complexity of meta learner depends heavily on the parameter quantity needing to be generated, thus increasing the difficulty of model training. Due to the ceiling effect of long-term dependence in sequence learning [174], the L2A approaches are difficult to generalize the case with slightly more support samples in a task.

4.4 Other Approaches

In addition to the aforementioned three mainstreams, i.e., augmentation (Section 4.1), metric learning (Section 4.2) and meta learning (Section 4.3), there have also some niche discriminative FSL approaches from other perspectives.

Multi-task learning [175] advocates learning multiple tasks synchronously by making the upstream embedding module implicitly or explicitly shared across tasks and the downstream task module specific in the hope of rendering the internal representation more generic. Follow this line, several multi-task learning based FSL approaches [176], [177], [178] have been proposed. In [176], a regularization penalty term is designed to force the parameters of different tasks to be similar. MetaGAN [177] introduced a task-conditioning GAN, which generates and discriminates fake samples and casts it as an auxiliary task to sharpen the decision boundary formed by other meta learning based FSL approaches. Z. Hu *et al.* [178] inserted an attribute prediction step before final class prediction and combined the attribute learning loss with the main task loss to jointly optimize the whole learner.

Self-supervised learning becomes a popular technique nowadays, especially in the field of vision [180], [181], [182], [183], [184], [185], [186], [187], [188], [189] and language [190], [191], [192], [193], which aims to learn semantically meaningful representations using only the inherent structural information contained by data itself instead of expensive human labels. Researchers leverage the structural information of data as its self-supervision to train their networks. For instance, in [181], [185], [186], an unlabeled image is shuffled into several patches and the permutation of patches is treated as the self-supervision, and thus the goal of self-supervised learning is to solve jigsaw puzzles. Besides, other self-supervised learning tasks include predicting image rotation angles [184], relative patch location [180], exemplar class of augmented samples [194], etc. Recently, there are also several works attempting to address FSL problems in the spirit of self-supervised learning [179], [195], [196], [197]. For example, as shown in Fig. 23, [179] combined the supervised loss formed by the off-the-shelf Prototypical Nets [54] and the

self-supervised losses formed by rotation tasks and jigsaw puzzle tasks to learn the feature representation. S2M2 [197] leveraged self-supervised tasks (i.e., rotation prediction and exemplar prediction) and Manifold Mixup [197] to regularize the feature manifold, which leads to an additional loss for general-purpose representation. Both [195], [196] focused on how to incorporate self-supervised learning into semi-supervised FSL tasks, which will be introduced in Section 5.1. In essence, these self-supervised FSL approaches construct auxiliary self-supervised tasks attached to the main FSL tasks, and thus they still belong to the scope of multi-task learning.

Inspired by transductive inference [198], Y. Liu *et al.* [199] assumed a transductive setting where all query samples in a task would arrive at once when testing. In this manner, Transductive Propagation Network (TPN) [199] was developed, which realized a label propagation for all unlabeled query samples using a graph model. Besides, several works proposed to fully utilize extra available data or prior knowledge to facilitate FSL. For instance, Z. Xu *et al.* [200] drew support from large-scale machine-labeled web images, and M. Bauer *et al.* [201] leveraged the concept information between different classes to build a probabilistic K -shot learning model. In addition, there also exist some approaches from other unique perspectives, such as feature replacement [40], LS-SVM-based model adaptation [41], Bilevel Programming [202], knowledge distillation [203], dense classification [204] and saliency-guided data hallucination [205], etc.

With the recent emergence of many FSL solutions, some researchers switched their focus from method development to further analysis of existing methods. In [206], the effect of the base network's depth to FSL model capability was analyzed via a series of consistent comparative experiments on several representative FSL approaches, such as Prototypical Nets [54], Matching Nets [49], Relation Net [55] and MAML [50]. In [207], the metric learning based FSL approaches were further explored and it was claimed that some simple feature pre-processing (e.g., mean-subtraction and L2-normalization) could bring performance improvement.

5 EXTENSIONAL TOPICS

This section elaborates several emerging extensional topics of FSL including Semi-supervised FSL (S-FSL), Unsupervised FSL (U-FSL), Cross-domain FSL (C-FSL), Generalized FSL (G-FSL) and Multimodal FSL (M-FSL). Their mathematical descriptions have been presented in Section 2.2. Based on a variety of practical application environments, the five topics replan the application scenarios and task requirements of FSL, and they are becoming the hot directions of FSL researches.

5.1 Semi-supervised Few Sample Learning

S-FSL postulates that the training set D_{trn} for an N -way K -shot task contains not only NK labeled support samples, but also some unlabeled samples that belong or not belong to the C task classes. Researchers are allowed to use the semi-supervised training set to build their FSL systems.

In [208], [209], [210], Prototypical Nets [54] were improved to cope with S-FSL via semi-supervised clustering. MetaGAN [177] that has been discussed in Section 4.4 was designed to be compatible with the S-FSL setting.

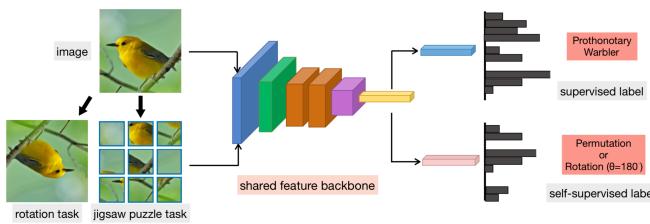


Fig. 23. Overview of self-supervised learning based FSL approach [179].

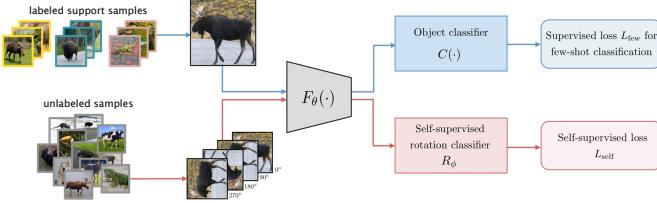


Fig. 24. Overview of self-supervised approach for S-FSL tasks [195].

In [195] and [196], the self-supervised learning paradigm was exploited to ingest information from unlabeled samples. In particular, as shown in Fig. 24, [195] constructed self-supervised tasks (i.e., rotation prediction and relative patch location) on unlabeled images and added this self-supervised loss into the main FSL task loss, which is very similar to the approach in [179] discussed in Section 4.4 except that it built the self-supervised tasks using labeled support samples (see Fig. 23). Differently, [196] proposed a self-training strategy for S-FSL that iterates between predicting pseudo labels for unlabeled data and finetuning FSL models using the pseudo-labeled data. Self-Jig [92] treated the labeled images as probe samples and the unlabeled images as gallery samples and then synthesized new images from them. Besides, several task-oriented S-FSL models have also been proposed recently, like SAMIE [211] for question-answer tasks and AffinityNet [212] for disease prediction tasks.

5.2 Unsupervised Few Sample Learning

U-FSL encourages a more general setting than vanilla FSL where the auxiliary set D_A is fully unsupervised. The goal is to pursue a relatively mild condition for performing FSL and weaken the prerequisite for building an FSL learner, since collecting an unlabeled auxiliary set belonging to non-task classes is more easy-to-implement than collecting a labeled dataset. For example, one can readily acquire a large number of unlabeled images via web crawlers in today's big data era.

In [151], the top layers in the base learner were pre-trained as low-density separator (LDS) using the unlabeled samples, and they are encouraged to capture a more generic representation space for downstream FSL tasks. CACTUs [213] adopted a two-stage strategy: synthesizing meta-train tasks on the unlabeled set by unsupervised representation learning methods (e.g., ACAI [214] and BiGAN [215]) and clustering algorithms, and then running classic MAML [50] or Prototypical Nets [54] on these synthetic tasks. Comparably, both UMTRA [215] and AAL [216] synthesized meta-train tasks through augmenting the unlabeled samples and treating the ancestor, on which augmentation is performed, and the corresponding augmented data as the congener samples, which is followed by the ready-made MAML [50] algorithm. The commonness between CACTUs [213], UMTRA [215] and AAL [216] is that they essentially focused on how to allocate pseudo labels to unlabeled samples such that the existing vanilla FSL models can work without modification.

5.3 Cross-domain Few Sample Learning

Under the vanilla FSL setting, it is assumed that the samples in auxiliary dataset D_A and T -specific dataset D_T all from the same data domain, as depicted in the top part of Fig. 25.

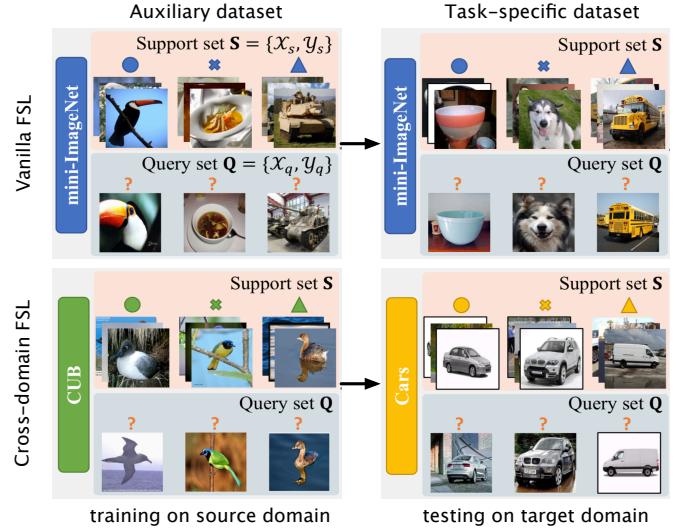


Fig. 25. Illustration for the task setting of Cross-domain FSL (C-FSL).

TABLE 6
Specific cross-domain forms studied by C-FSL approaches

Approaches	Cross-domain Form ($D_A \rightarrow D_T$)
[219]	SVHN [225] 0-4 → MNIST [60] 5-9, ImageNet [23] → UCF-101 [226]
[220]	Omniglot [37] → EMNIST [90]
	<i>Digit dataset:</i> MNIST [60] → USPS [227], MNIST → SVHN [225] USPS → MNIST, SVHN → USPS
[221]	<i>Office dataset</i> [228]: Amazon → DSLR, Amazon → Webcam, DSLR → Webcam, Webcam → DSLR
[222]	<i>Character dataset:</i> Omniglot → Omniglot-M, Omniglot-M → Omniglot
	<i>Office-Home dataset</i> [229]: Clipart → Product, Product → Clipart
[164]	<i>Office-Home dataset</i> [229]: Clipart → Real World, Real World → Clipart
[223]	miniImageNet [49] → CUB [82] / Cars [230] / Places [231] / Plantae [232]
[224]	MNIST [60] → Cifar-10 [81], Cifar-10 → MNIST
[206]	miniImageNet [49] → CUB [82]
[119]	miniImageNet [49] → CUB [82], Kinetics-CMN [233] → Jester
[122]	miniImageNet [49] → CUB [82]

However, when the FSL task to be handled is from a novel domain for which no relevant auxiliary samples are available, we have to leverage some cross-domain samples as the auxiliary data, as shown in the bottom part of Fig. 25. The domain shift between auxiliary dataset and task-specific dataset poses a higher challenge for FSL approaches.

C-FSL is naturally highly related to domain adaptation (DA) [217], which is a classic direction in the field of machine learning. Although there exists individual work [218] addressing DA with a few samples, its task setting is different from C-FSL: the label space in DA is shared between source and target domain, whereas that in C-FSL tasks is disjoint between auxiliary dataset and task-specific dataset. Recently, several approaches were proposed to tackle C-FSL problems from various perspectives, such as adversarial training [219], [220], [221], [222], feature transformation [164], [223], domain alignment [224], domain-specific finetuning [206], feature composition [119] and ensemble methods [122], etc. To facilitate follow-up C-FSL related researches, we summarize in Table 6 the specific cross-domain forms used by them.

5.4 Generalized Few Sample Learning

Vanilla FSL setting could easily lead to catastrophic forgetting issue [234], that is, most FSL models were trained to make inference for pre-defined classes of a novel task, but can not be continuously applied to the previous classes in the auxiliary set. However, in many applications where class concepts and samples arrive in a dynamic manner, learning systems are often faced with an extreme imbalance of training data among classes, which means some classes are provided with sufficient training samples while some have only a few. In this context, it is crucial and desirable to have the incremental learning ability for novel task classes with limited data while at the same time does not forget previous non-task classes. Thus, the focus of G-FSL is to enable FSL models to jointly handle all classes in both D_A and D_T .

Several augmentation based FSL approaches mentioned in Section 4.1 including SH [85], Hallucinator [86], CP-ANN [87] and IDMe-Net [92] are naturally applicable to both FSL and G-FSL settings since their learning processes were divided into two independent stages: first augmenting training samples for the sparse task classes and then training the models using the combination of raw and augmented samples. GcGPN [235] extended vanilla Prototypical Nets [54] to G-FSL setting using a GCN [167], which models the relationship between all novel and existing classes by casting the classes as nodes and the inter-class dependencies as edges. Dynamic Nets [152], Acts2Params [153], DAE [157] and AAN [236] adhered to a common principle of generating new weights for novel classes incrementally and combining them into the weights of existing classes to form a joint decision maker. In [237], a notion of class adapting principal directions is introduced to enable efficient and discriminative embeddings for images from both novel and existing classes. CADA-VAE [238] developed a Variational Auto-Encoder to create latent space features for novel classes and then trained the final predictor over all classes. The L2A approach, FEAT [164], is also suitable for G-FSL setting since its set-to-set functions are class-agnostic and its embedding adaptation can operate on both novel and existing classes.

5.5 Multimodal Few Sample Learning

Different from vanilla FSL that only contains a task modality, M-FSL involves information or data from additional modality. According to the role of additional modality, M-FSL setting can be further subdivided into two cases, as shown in Fig. 26.

Multimodal Matching. Vanilla FSL seeks for a mapping from the task modality to hard class label space, whereas the multimodal matching of M-FSL aims to learn the mapping from one modality to another modality [239], [240], [241]. For example, given a few image-sentence training pairs, FSL learners are required to determine the sentence correctly describing the query image [239], [240], or given a handful of speech-images training pairs, FSL learners need to find a correct visual image that contains the word spoken in the query speech [241]. These FSL based multimodal matching settings are meaningful, especially for robotic applications.

Multimodal Fusion. It allows FSL learners to use extra information from additional modalities to help the learning in task modality. Several recent works [116], [155], [238], [242], [243], [244], [245], [246], [247], [248] enhanced FSL

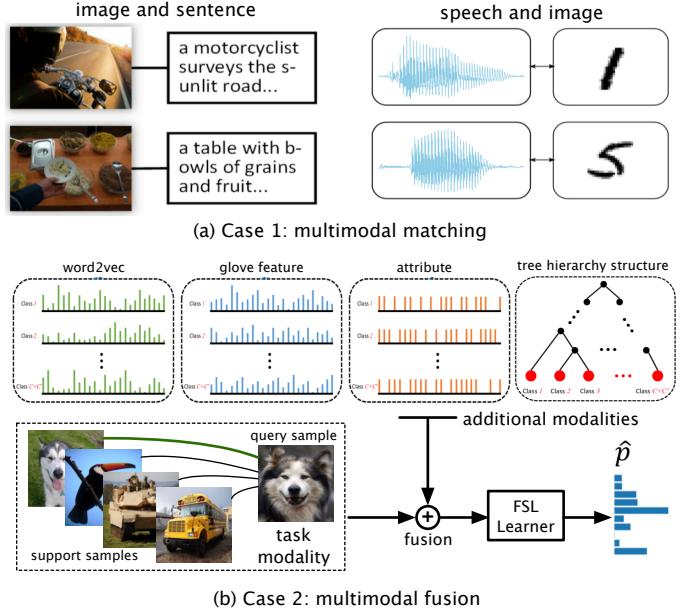


Fig. 26. Illustration for the two cases of Multimodal FSL (M-FSL).

model capability by fusing various multimodal information, which include word2vec [238], [244], text captions [155], [242], [243], [247], attributes [238], [244], [249], glove features [244], word embeddings [116], [245], tree hierarchy structure among classes [244], [248] and cross-style datasets [246], etc. These additional modalities bring more prior knowledge into FSL, providing a remedy for the lack of training samples.

6 APPLICATIONS

Since the ubiquitous demand of machine learning systems for large-scale training samples and the vigorous advances of FSL studies in recent years, the methods and ideas of FSL are being widely applied to various research areas such as computer vision, natural language processing, audio and speech, reinforcement learning and robotic, and data analysis, etc. Table 9 summarizes the fields and subfields of FSL applications as well as their representative publications.

TABLE 7
Statistics of popular FSL benchmark datasets for image classification

Dataset	# Images	# Trn/Val/Tst cls.	Content
Omniglot [37]	129,840	4,800/-/1,692	characters
miniImageNet [49]	60,000	64/16/20	common objects
tieredImageNet [209]	779,165	351/97/160	common objects
CUB [127]	11,788	100/50/50	birds
Stanford Dogs [130]	20,580	70/20/30	dogs
Stanford Cars [130]	16,185	130/17/49	cars
Caltech-256 [140]	30,607	150/56/50	common objects
Oxford-102 [126]	8,189	82/-/20	flowers
FC100 [114]	60,000	60/20/20	common objects
CIFAR-FS [159]	60,000	64/26/20	common objects
Visual Genome [245]	~108,000	1,211/-/829	common objects
SUN397 [249]	108,754	197/-/200	scenes
ImageNet1K [85]	~1,000,000	389/-/611	common objects

TABLE 8

Image classification accuracy (%) on *miniImageNet* [49] of existing FSL approaches. Results are cited from their original articles, which were reported by averaging hundreds of random meta-test C -way K -shot tasks with 95% confidence intervals ("—" not available)

Approaches	5-way 1-shot	5-way 5-shot	Approaches	5-way 1-shot	5-way 5-shot
Matching Nets [49]	43.56 ± 0.84	55.31 ± 0.73	Resnet PN [208]	54.05 ± 0.47	70.92 ± 0.66
Meta-Learner LSTM [51]	43.44 ± 0.77	60.60 ± 0.71	MetaHebb [163]	56.84 ± 0.52	71.00 ± 0.34
MAML [50]	48.70 ± 1.84	63.11 ± 0.92	STANet [250]	58.35 ± 0.57	71.07 ± 0.39
MACO [127]	41.09 ± -	58.32 ± -	CSNs [162]	56.88 ± 0.62	71.94 ± 0.57
Gauss (MAP pr.) HMC [201]	50.00 ± 0.50	64.30 ± 0.60	SalNet [205]	57.45 ± 0.88	72.01 ± 0.67
Meta-SGD [137]	50.47 ± 1.87	64.03 ± 0.94	Dynamic Nets [152]	56.20 ± 0.86	72.81 ± 0.62
Reptile [139]	48.21 ± 0.69	66.00 ± 0.62	Dual TriNet [74]	58.12 ± 1.37	76.92 ± 0.69
MetaNet [50]	49.21 ± 0.96	-	Act2Params [153]	59.60 ± 0.41	73.74 ± 0.19
LLAMA [141]	49.40 ± 1.83	-	TADAM [114]	58.50 ± 0.30	76.70 ± 0.30
Prototypical Nets [54]	49.42 ± 0.78	68.20 ± 0.66	Deep Comparison Net [128]	62.88 ± 0.83	75.84 ± 0.65
IMP [210]	49.60 ± 0.80	68.10 ± 0.80	IDeMe-Net [92]	59.14 ± 0.86	74.63 ± 0.74
GNN [133]	50.33 ± 0.36	66.41 ± 0.63	K-tuple Nets [120]	58.30 ± 0.84	72.37 ± 0.63
Triplet Ranking Nets [98]	50.58 ± -	-	Self-Jig [92]	58.80 ± 1.36	76.71 ± 0.72
mAP-Nets [111]	50.32 ± 0.80	63.94 ± 0.72	CAML [146]	59.23 ± 0.99	72.35 ± 0.71
Relation Net [55]	50.44 ± 0.82	65.32 ± 0.70	CFA [119]	58.50 ± 0.80	76.60 ± 0.60
Cross-Modulation Nets [123]	50.94 ± 0.61	66.65 ± 0.67	SoSN [131]	59.22 ± 0.91	73.24 ± 0.69
Hyper-Represent [202]	50.54 ± 0.85	64.53 ± 0.68	DAE [157]	61.07 ± 0.15	76.75 ± 0.11
CovaMNet [129]	51.19 ± 0.76	67.65 ± 0.63	LEO [145]	61.76 ± 0.08	77.59 ± 0.12
TAML [144]	51.73 ± 1.88	66.05 ± 0.85	AAM [118]	62.24 ± 0.20	77.24 ± 0.15
Large Margin [135]	51.41 ± 0.68	67.81 ± 0.64	MTL [149]	61.20 ± 1.80	75.50 ± 0.80
SARN [132]	51.62 ± 0.31	66.16 ± 0.51	EGNN [134]	-	76.37 ± -
MT-net [138]	51.70 ± 1.84	-	Principal Characteristic Nets [121]	63.29 ± 0.76	77.08 ± 0.68
MM-Net [124]	53.37 ± 0.48	66.97 ± 0.35	AM3 [116]	65.30 ± 0.49	78.10 ± 0.36
MetaGAN [177]	52.71 ± 0.64	68.63 ± 0.67	DC [204]	62.53 ± 0.19	78.95 ± 0.13
VERSA [158]	53.40 ± 1.82	67.37 ± 0.86	CC+rot [195]	62.93 ± 0.45	79.87 ± 0.33
BMAML [143]	53.80 ± 1.46	-	MetaOptNet [160]	64.09 ± 0.62	80.00 ± 0.45
SNAIL [171]	55.71 ± 0.99	68.88 ± 0.92	CTM [136]	64.12 ± 0.82	80.51 ± 0.13
DA-PN [224]	50.56 ± 0.85	69.62 ± 0.76	LGM-Net [56]	69.13 ± 0.35	71.18 ± 0.68
R2-D2 [159]	51.90 ± 0.20	68.70 ± 0.20	Diversity with Cooperation [122]	63.73 ± 0.62	81.19 ± 0.43
TPN [199]	55.51 ± -	69.86 ± -	FEAT [164]	66.78 ± -	82.05 ± -
SRPN [99]	55.20 ± -	69.60 ± -	SimpleShot [207]	64.29 ± 0.20	81.50 ± 0.14
Δ -encoder [88]	59.90 ± -	69.70 ± -	S2M2 [197]	64.93 ± 0.18	83.18 ± 0.11*
DN4 [130]	51.24 ± 0.74	71.02 ± 0.64	LST [196]	70.10 ± 1.90*	78.70 ± 0.80

TABLE 9
Summary of FSL applications in various fields and their representative publications

Fields	Subfields & References	
Image	image classification	general image classification (see Table 1, 2, 3, 4, 5, 6, 8, Fig. 13, 16), multi-label classification [251], fine-grained recognition [119], [129], [130], [243], [247], [252], [253], [254], [255], hyperspectral image classification [256], [257], 3D object/model classification [148], [251]
	image segmentation	semantic segmentation [258], [259], [260], [261], [262], [263], [264], [265], [266], instance segmentation [267], [267], [268], texture segmentation [269], [270], medical/biological image segmentation [271], [272], [273], [274]
	object detection	general objects [275], [276], [277], [278], air vehicles [279], RGB-D objects [280]
Computer Vision	other applications	image generation [37], [48], [69], [252], [281], [282], [283], [284], [285], image retrieval [111], [286], gaze estimation [287], depth estimation [288], localization [289], scene graph prediction [290], image-based person re-identification [291], [292], image colorization [293], color constancy [294]
	video classification	general video classification [176], [233], [295], [296], gesture recognition [297], [298], action recognition [119], [299], [300], [301], [302], [303], [304]
Video	video detection	action localization [305], [306], activity detection [307]
	other applications	video prediction [156], [308], video object segmentation [309], [310], semantic indexing [311], video retargeting [312], video generation [313], video-based person re-identification [314], object tracking [315], motion capture [316]
Natural Language Processing	text classification [317], [318], [319], [320], [321], [322], [323], [324], [325], [326], [327], [328], dialogue system [329], [330], [331], relation learning and knowledge graphs [332], [333], [334], [335], [336], word representation learning [337], [338], [339], [340], [341], named entity recognition [342], [343], [344], word prediction [49], [162], [163], natural language generation [345], [346], [347], [348], information extraction [211], machine translation [100], charge prediction [178], sequence labeling [349]	
Audio&Speech	audio/speech/sound classification [350], [351], [352], [353], [354], [355], text-to-speech [356], [357], [358], [359], acoustic/sound event detection [360], [361], [362], speech generation [350], [363], keyword/command recognition [364], keyword spotting [365], human-fall detection [366], speaker recognition [367]	
Reinforcement Learning&Robotic	imitation learning [368], [369], [370], [371], [372], [373], [374], [375], [376], locomotion [50], [171], [377], policy learning [246], [306], visual navigation [50], [137], [144], [171], [378], robot manipulation [378], [379], multi-armed bandits [171], tabular MDPs [171]	
Data Analysis	data regression [50], [137], [138], [139], [141], [142], [143], [145], [156], [169], [246], anomaly/error detection [380], [381], [382]	
Cross-Field	image captioning [383], visual question answering [383], [384]	
Other Applications	disease prediction [212], [385], [386], [387], [388], [389], [390], biometrical recognition (e.g., palmprint [391], ear [205]), drug discovery [392], spectrum classification [393], precision agriculture [394], internet security [395], mobile sensing [396]	

Computer Vision. Thanks to the intuitiveness and intelligibility of visual data, computer vision has always been the main testbed for machine learning algorithms, and it is no exception to FSL. From the earliest Congealing model [29] to today’s meta learning approaches, visual tasks have always acted as the touchstone of FSL approaches, especially the few sample based (or few-shot) image classification tasks. In Table 7, we enumerate several popular FSL benchmark datasets for image classification and summarize their statistics. Two most commonly used benchmarks are Omniglot [37] and *miniImageNet* [49]. Due to the simplicity of grayscale characters images and the sufficient meta-train classes, many FSL approaches have achieved good performance on Omniglot that is close to saturation. As a result, researchers are inclined to utilize *miniImageNet* to evaluate the performance of FSL approaches. For better reference to follow-up studies, we summarize in Table 8 the performance of all FSL approaches that have reported their results on *miniImageNet*. We can observe that, within only three years from 2016 to 2019, the 5-way 1/5-shot accuracy increased by more than 20%, which indicates the rapid development of FSL researches. Besides, FSL has been incorporated into image segmentation [397], object detection [398] and other image-based vision tasks. At the level of video data, FSL also has many rising applications in video classification [399], video detection [400], video object segmentation [401], etc. More FSL applications in the vision domain can be found in the first part of Table 9.

Natural Language Processing. It is the second largest field of FSL applications. One common FSL application in natural language processing is text classification [402], which seeks to utilize a few documents or words to infer document labels. In addition, the FSL regime was also brought into fundamental research topics of natural language processing, such as word representation learning [403], relation learning and knowledge graphs [404]. The second part of Table 9 details more FSL applications in natural language processing.

Audio and Speech. Acoustic data is a more complex data form, and generally the large-scale collection and annotation for them are more difficult than that for images or texts, which leads to a more urgent need for FSL approaches. At present, FSL has been used to address many acoustic tasks covering from the basic audio classification and keyword recognition to the challenging text-to-speech and speech generation. The third part of Table 9 summarizes existing FSL applications and corresponding FSL references.

Reinforcement Learning and Robotic. An ideal robotic system should possess the ability of learning novel tasks with a few demonstrations and without long task-specific training time for a task, however, a new situation could make robots vulnerable to the dilemma of limited observation samples, which makes FSL an indispensable skill for future advanced robotic systems. As FSL approaches grew in popularity, many researchers have reconsidered the applications of reinforcement learning and robotic [405] under the regime of FSL, which include imitation learning [406], visual navigation [407] and policy learning [408], etc. More related applications are presented in the fourth part of Table 9.

Data Analysis. As is well known, effectively analyzing data and mining underlying rules in data via sparse training data is a goal tirelessly pursued by data science researchers.

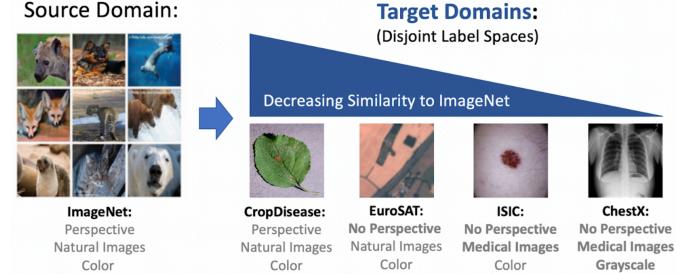


Fig. 27. Illustration for the cross-domain few-shot learning competition (<https://www.learning-with-limited-labels.com/challenge>).

Fortunately, FSL is gradually being applied to some classical data analysis applications like data regression and anomaly detection [409], as described in the fifth part of Table 9.

Cross-Field Applications. Recently, FSL has been integrated into two popular cross-field applications, i.e., image captioning [383] and visual question answering [383], [384]. Given only a few image-text training pairs, the former tries to generate a proper textual description for an image, while the latter seeks to output an accurate natural language answer to a textual question about an image.

Other Applications. Besides the above several common application fields of machine learning, FSL has also been introduced into other professional areas, such as medicine, chemometrics, agriculture, sensors and internet security, etc. Please refer to the last part of Table 9 for more details.

Open Competitions. With the growing attention towards FSL, several related competitions are emerging. To our best knowledge, *Few-Shot Verb Image Classification* (<http://www.lsfsl.net/cl/>) that was published in a workshop of ICCV 2019 is the first FSL competition, which focused on large-scale verb image classification and proposed a high-quality few-shot verb image dataset. Recently, in the Visual Learning with Limited Labels Workshop of CVPR 2020, a more challenging competition, *Cross-Domain Few-Shot Learning Challenge* (<https://www.learning-with-limited-labels.com/challenge>), was proposed, which keeps aligned with the C-FSL tasks discussed in Section 5.3. As depicted in Fig. 27, this competition requires participants to train FSL models on ImageNet but perform evaluation on the other four datasets from varying domains, such as plant disease images, satellite images, dermoscopic images of skin lesions, and X-ray images. This competition includes two main tracks that use or not use unlabeled images from the target domain for training.

7 FUTURE DIRECTIONS

Though recent years have witnessed considerable progress of FSL in both methodology and applications, challenges still exist due to the intrinsic difficulty from sparse samples. In this section, we suggest four future directions of FSL.

Robustness. Most of the current FSL studies are based on an ideal data hypothesis, but it is hard to hold true for all practical scenes. In many realistic applications, one may be faced with uncertain disturbances that destroy the ideal setting of FSL. For instance, the few training data may suffer from outlier interference (e.g., noisy samples or label-wrong data) [410] due to instrumental malfunction or

perfunctory errors. It raises the question of whether existing FSL models can effectively alleviate the influence from such outliers and still maintain an acceptable generalization. In addition, the possible domain shift between auxiliary data and task-specific data as described in Section 5.3 is another kind of disturbance to the ideal setting of FSL. Therefore, improving the robustness of FSL models against various potential disturbance factors is substantially meaningful.

Universality. The universality mentioned here is twofold. The first is the model-level generality and scalability of FSL approaches. For now, most of FSL approaches are excessively designed for the specific benchmark tasks and datasets, weakening their applicability to other more general tasks. An ideal FSL framework should be able to deal with various learning tasks with different data complexity and diverse data forms. The second is the application-level versatility and flexibility of FSL approaches. The majority of current FSL studies focus on the plain application scenario with small-scale task classes and large-scale labeled auxiliary data. However, real-world problems may bring more complex application scenarios, such as large-scale task classes, long-tail phenomenon of data distribution [411], dynamicity of task classes, unavailability of labeled auxiliary data, and even a mixture of these scenarios. They raise higher requirements and challenges for the universality of FSL approaches.

Interpretability. The surge and success of FSL in recent years mainly lie in deep learning technology, which is often criticized for its lack of interpretability. Model interpretability is a key issue for deep learning [412], [413]. We believe that the impressive few sample learning ability of humans benefit from many aspects including the rational use of empirical knowledge and the ingenious exploration of underlying knowledge behind task data (e.g., compositional relationship [37], [63], structural correspondence between data components [414], etc). Therefore, how to capitalize on the fusion of external prior knowledge and internal data knowledge to enhance the interpretability of FSL models could be a future research direction.

Theoretical System. As analyzed in Section 1, the fundamental difficulty caused by sparse training samples is that the search space of learning function f is very huge due to the lack of effective function regularization formed by training samples. If we re-look current FSL approaches from this theoretical view of point, it can be found that, in essence, all FSL solutions are to realize function regularization through specific technologies. For instance, the augmentation based FSL approaches reach this goal by directly increasing the training samples, while the meta learning approaches suggest introducing other irrelevant learning tasks to regularize the learning function across tasks. Thus, building a systematic theoretical system for FSL from the perspective of regularizing learning function space under sparse training samples could bring new inspiration to FSL researchers.

8 CONCLUSIONS

Enabling learning systems to learn from very few samples is crucial for the further development of machine learning and artificial intelligence. This article conducts a comprehensive survey on few sample learning (FSL). In particular, the evolution history and current advances of FSL are reviewed,

and all FSL approaches are grouped via a succinct and understandable taxonomy. An in-depth analysis is made to shed light into the underlying development relationship between mainstream meta learning based FSL approaches. Also, several emerging extensional research topics of FSL, existing FSL applications in various fields, current benchmark datasets and performance, together with several potential research directions are systematically summarized. This survey is expected to promote the grasp of FSL related knowledge and the collaborative development of FSL research area.

ACKNOWLEDGMENTS

The authors would like to thank the pioneer researchers in few sample learning and other related fields. This work was funded in part by the National Natural Science Foundation of China under Grant 61876095 and Grant 61751308, in part by the Beijing Natural Science Foundation, and in part by the Beijing Academy of Artificial Intelligence (BAAI).

REFERENCES

- [1] B. Landau, L. B. Smith, and S. S. Jones, "The importance of shape in early lexical learning," *Cognitive Development*, vol. 3, no. 3, pp. 299–321, 1988.
- [2] E. M. Markman, *Categorization and naming in children: Problems of induction*. MIT Press, 1989.
- [3] F. Xu and J. B. Tenenbaum, "Word learning as bayesian inference." *Psychological Review*, vol. 114, no. 2, pp. 245–272, 2007.
- [4] I. Biederman, "Recognition-by-components: a theory of human image understanding." *Psychological Review*, vol. 94, no. 2, pp. 115–147, 1987.
- [5] S. Carey and E. Bartlett, "Acquiring a single new word." *Papers and Reports on Child Language Development*, vol. 15, pp. 17–29, 1978.
- [6] E. V. Clark, *First language acquisition*. Cambridge University Press, 2009.
- [7] N. P. Rougier, D. C. Noelle, T. S. Braver, J. D. Cohen, and R. C. O'Reilly, "Prefrontal cortex and flexible cognitive control: Rules without symbols," *Proc. Nation. Academy Sci.*, vol. 102, no. 20, pp. 7338–7343, 2005.
- [8] T. S. Braver, J. L. Paxton, H. S. Locke, and D. M. Barch, "Flexible neural mechanisms of cognitive control within human prefrontal cortex," *Proc. Nation. Academy Sci.*, vol. 106, no. 18, pp. 7351–7356, 2009.
- [9] A. A. Kehagia, G. K. Murray, and T. W. Robbins, "Learning and cognitive flexibility: frontostriatal function and monoaminergic modulation," *Current Opinion in Neurobiology*, vol. 20, no. 2, pp. 199–204, 2010.
- [10] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [11] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2012, pp. 1097–1105.
- [12] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, pp. 1–9.
- [13] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 770–778.
- [14] T. Mikolov, M. Karafiat, L. Burget, J. Černocký, and S. Khudanpur, "Recurrent neural network based language model," in *Proc. 11th Annu. Conf. Int. Speech Commun. Assoc. (INTERSPEECH)*, 2010.
- [15] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2014, pp. 3104–3112.
- [16] G. Hinton, L. Deng, D. Yu, G. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, B. Kingsbury *et al.*, "Deep neural networks for acoustic modeling in speech recognition," *IEEE Signal Process. Mag.*, vol. 29, no. 6, pp. 82–97, 2012.

- [17] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [18] T. Gebru, J. Krause, Y. Wang, D. Chen, J. Deng, E. L. Aiden, and L. Fei-Fei, "Using deep learning and google street view to estimate the demographic makeup of neighborhoods across the united states," *Proc. Nation. Academy Sci.*, vol. 114, no. 50, pp. 13 108–13 113, 2017.
- [19] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun, "Dermatologist-level classification of skin cancer with deep neural networks," *Nature*, vol. 542, no. 7639, pp. 115–118, 2017.
- [20] S. Ghosal, D. Blystone, A. K. Singh, B. Ganapathysubramanian, A. Singh, and S. Sarkar, "An explainable deep machine vision framework for plant stress phenotyping," *Proc. Nation. Academy Sci.*, vol. 115, no. 18, pp. 4613–4618, 2018.
- [21] M. S. Norouzzadeh, A. Nguyen, M. Kosmala, A. Swanson, M. S. Palmer, C. Packer, and J. Clune, "Automatically identifying, counting, and describing wild animals in camera-trap images with deep learning," *Proc. Nation. Academy Sci.*, vol. 115, no. 25, pp. E5716–E5725, 2018.
- [22] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [23] 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," *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, 2015.
- [24] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *Int. J. Comput. Vis.*, vol. 88, no. 2, pp. 303–338, 2010.
- [25] S. Legg and M. Hutter, "Universal intelligence: A definition of machine intelligence," *Minds and Machines*, vol. 17, no. 4, pp. 391–444, 2007.
- [26] W. Wang, V. W. Zheng, H. Yu, and C. Miao, "A survey of zero-shot learning: Settings, methods, and applications," *ACM Trans. Intell. Syst. Technology*, vol. 10, no. 2, p. 13, 2019.
- [27] J. Shu, Z. Xu, and D. Meng, "Small sample learning in big data era," *arXiv preprint arXiv:1808.04572*, 2018.
- [28] Y. Wang, Q. Yao, J. T. Kwok, and N. L. M., "Generalizing from a few examples: a survey on few-shot learning," *arXiv preprint arXiv:1904.05046*, 2019.
- [29] E. G. Miller, N. E. Matsakis, and P. A. Viola, "Learning from one example through shared densities on transforms," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2000, pp. 464–471.
- [30] G. Koch, R. Zemel, and R. Salakhutdinov, "Siamese neural networks for one-shot image recognition," in *Proc. Int. Conf. Mach. Learn. (ICML) Deep Learn. Workshop*, vol. 2, 2015.
- [31] L. Fei-Fei, R. Fergus, and P. Perona, "A bayesian approach to unsupervised one-shot learning of object categories," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2003, pp. 1134–1141.
- [32] L. Fei-Fei, R. Fergus, and P. Perona, "Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR) Workshop*, 2004.
- [33] ——, "One-shot learning of object categories," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 4, pp. 594–611, 2006.
- [34] B. Lake, R. Salakhutdinov, J. Gross, and J. Tenenbaum, "One shot learning of simple visual concepts," in *Proc. Annu. Meet. Cognit. Sci. Soc. (CogSci)*, vol. 33, no. 33, 2011.
- [35] B. Lake, R. Salakhutdinov, and J. Tenenbaum, "Concept learning as motor program induction: A large-scale empirical study," in *Proc. Annu. Meet. Cognit. Sci. Soc. (CogSci)*, vol. 34, no. 34, 2012.
- [36] B. M. Lake, R. R. Salakhutdinov, and J. Tenenbaum, "One-shot learning by inverting a compositional causal process," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2013, pp. 2526–2534.
- [37] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum, "Human-level concept learning through probabilistic program induction," *Science*, vol. 350, no. 6266, pp. 1332–1338, 2015.
- [38] M. Fink, "Object classification from a single example utilizing class relevance metrics," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2005, pp. 449–456.
- [39] L. Wolf and I. Martin, "Robust boosting for learning from few examples," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2005, pp. 359–364.
- [40] E. Bart and S. Ullman, "Cross-generalization: Learning novel classes from a single example by feature replacement," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2005, pp. 672–679.
- [41] T. Tommasi and B. Caputo, "The more you know, the less you learn: from knowledge transfer to one-shot learning of object categories," in *Proc. British Mach. Vis. Conf. (BMVC)*, 2009.
- [42] X. Yu and Y. Aloimonos, "Attribute-based transfer learning for object categorization with zero/one training example," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*. Springer, 2010, pp. 127–140.
- [43] K. D. Tang, M. F. Tappen, R. Sukthankar, and C. H. Lampert, "Optimizing one-shot recognition with micro-set learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2010, pp. 3027–3034.
- [44] M. A. Tanner and W. H. Wong, "The calculation of posterior distributions by data augmentation," *J. American Statistical Assoc.*, vol. 82, no. 398, pp. 528–540, 1987.
- [45] E. P. Xing, M. I. Jordan, S. J. Russell, and A. Y. Ng, "Distance metric learning with application to clustering with side-information," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2003, pp. 521–528.
- [46] R. Vilalta and Y. Drissi, "A perspective view and survey of meta-learning," *Artificial Intell. Review*, vol. 18, no. 2, pp. 77–95, 2002.
- [47] H. Edwards and A. Storkey, "Towards a neural statistician," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2016.
- [48] D. Rezende, I. Danihelka, K. Gregor, D. Wierstra *et al.*, "One-shot generalization in deep generative models," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2016, pp. 1521–1529.
- [49] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, "Matching networks for one shot learning," in *Proc. Advances Neural Inf. Process. Syst. (NIPS)*, 2016, pp. 3630–3638.
- [50] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2017, pp. 1126–1135.
- [51] S. Ravi and H. Larochelle, "Optimization as a model for few-shot learning," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2017.
- [52] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap, "Meta-learning with memory-augmented neural networks," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2016, pp. 1842–1850.
- [53] T. Munkhdalai and H. Yu, "Meta networks," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2017, pp. 2554–2563.
- [54] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," in *Proc. Advances Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 4077–4087.
- [55] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales, "Learning to compare: Relation network for few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 1199–1208.
- [56] H. Li, W. Dong, X. Mei, C. Ma, F. Huang, and B.-G. Hu, "Lgm-net: Learning to generate matching networks for few-shot learning," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2019, pp. 3825–3834.
- [57] P. J. Grother, "Nist special database 19-hand-printed forms and characters database," National Institute of Standards and Technology, Tech. Rep., 1995.
- [58] M. Weber, M. Welling, and P. Perona, "Unsupervised learning of models for recognition," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*. Springer, 2000, pp. 18–32.
- [59] R. Salakhutdinov, J. Tenenbaum, and A. Torralba, "One-shot learning with a hierarchical nonparametric bayesian model," in *Proc. Int. Conf. Mach. Learn. (ICML) Workshop on Unsupervised and Transfer Learning*, 2012, pp. 195–206.
- [60] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [61] F. Fleuret and G. Blanchard, "Pattern recognition from one example by chopping," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2006, pp. 371–378.
- [62] S. A. Nene, S. K. Nayar, and H. Murase, "Columbia object image library (coil-100)," CUCS-006-96, Columbia University, Tech. Rep., 1996.
- [63] A. Wong and A. L. Yuille, "One shot learning via compositions of meaningful patches," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2015, pp. 1197–1205.
- [64] J. J. Hull, "A database for handwritten text recognition research," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 16, no. 5, pp. 550–554, 1994.
- [65] H. Edwards and A. Storkey, "Towards a neural statistician," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2017.

- [66] L. Wolf, T. Hassner, and I. Maoz, "Face recognition in unconstrained videos with matched background similarity," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2011, pp. 529–534.
- [67] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," in *Proc. British Mach. Vis. Conf. (BMVC)*, 2014.
- [68] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*. Springer, 2014, pp. 818–833.
- [69] A. Antoniou, A. Storkey, and H. Edwards, "Data augmentation generative adversarial networks," in *Proc. Int. Conf. Learn. Represent. (ICLR) Workshop*, 2018.
- [70] C. H. Lampert, H. Nickisch, and S. Harmeling, "Attribute-based classification for zero-shot visual object categorization," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 3, pp. 453–465, 2013.
- [71] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2013, pp. 3111–3119.
- [72] R. Kwitt, S. Hegenbart, and M. Niethammer, "One-shot learning of scene locations via feature trajectory transfer," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 78–86.
- [73] M. Dixit, R. Kwitt, M. Niethammer, and N. Vasconcelos, "Agg: Attribute-guided augmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 7455–7463.
- [74] Z. Chen, Y. Fu, Y. Zhang, Y.-G. Jiang, X. Xue, and L. Sigal, "Semantic feature augmentation in few-shot learning," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018.
- [75] ——, "Multi-level semantic feature augmentation for one-shot learning," *IEEE Trans. Image Process.*, vol. 28, no. 9, pp. 4594–4605, 2019.
- [76] J. Lu, J. Li, Z. Yan, F. Mei, and C. Zhang, "Attribute-based synthetic network (abs-net): Learning more from pseudo feature representations," *Pattern Recognit.*, vol. 80, pp. 129–142, 2018.
- [77] P.-Y. Laffont, Z. Ren, X. Tao, C. Qian, and J. Hays, "Transient attributes for high-level understanding and editing of outdoor scenes," *ACM Trans. Graphics*, vol. 33, no. 4, pp. 1–11, 2014.
- [78] G. Patterson, C. Xu, H. Su, and J. Hays, "The sun attribute database: Beyond categories for deeper scene understanding," *Int. J. Comput. Vis.*, vol. 108, no. 1-2, pp. 59–81, 2014.
- [79] S. Song, S. P. Lichtenberg, and J. Xiao, "Sun rgb-d: A rgb-d scene understanding benchmark suite," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, pp. 567–576.
- [80] C. H. Lampert, H. Nickisch, and S. Harmeling, "Learning to detect unseen object classes by between-class attribute transfer," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2009, pp. 951–958.
- [81] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," Citeseer, Princeton, NJ, USA, Tech. Rep., 2009.
- [82] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "The caltech-ucsd birds-200-2011 dataset," California Institute of Technology, Tech. Rep., 2011.
- [83] G. Griffin, A. Holub, and P. Perona, "Caltech-256 object category dataset," 2007.
- [84] R. Fergus, P. Perona, and A. Zisserman, "Object class recognition by unsupervised scale-invariant learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, vol. 2, 2003, pp. II–II.
- [85] B. Hariharan and R. Girshick, "Low-shot visual recognition by shrinking and hallucinating features," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017, pp. 3018–3027.
- [86] Y.-X. Wang, R. Girshick, M. Hebert, and B. Hariharan, "Low-shot learning from imaginary data," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 7278–7286.
- [87] H. Gao, Z. Shou, A. Zareian, H. Zhang, and S.-F. Chang, "Low-shot learning via covariance-preserving adversarial augmentation networks," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 975–985.
- [88] E. Schwartz, L. Karlinsky, J. Shtok, S. Harary, M. Marder, A. Kumar, R. Feris, R. Giryes, and A. Bronstein, "Delta-encoder: an effective sample synthesis method for few-shot object recognition," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 2845–2855.
- [89] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth, "Describing objects by their attributes," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2009, pp. 1778–1785.
- [90] G. Cohen, S. Afshar, J. Tapson, and A. Van Schaik, "Emnist: Extending mnist to handwritten letters," in *Proc. Int. Joint Conf. Neural Net. (IJCNN)*, 2017, pp. 2921–2926.
- [91] O. M. Parkhi, A. Vedaldi, and A. Zisserman, "Deep face recognition," in *Proc. British Mach. Vis. Conf. (BMVC)*, 2015.
- [92] Z. Chen, Y. Fu, Y.-X. Wang, L. Ma, W. Liu, and M. Hebert, "Image deformation meta-networks for one-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 8680–8689.
- [93] M. Rosen-Zvi, C. Chemudugunta, T. Griffiths, P. Smyth, and M. Steyvers, "Learning author-topic models from text corpora," *ACM Trans. Inf. Syst.*, vol. 28, no. 1, pp. 1–38, 2010.
- [94] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2014, pp. 2672–2680.
- [95] T. Hertz, A. B. Hillel, and D. Weinshall, "Learning a kernel function for classification with small training samples," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2006, pp. 401–408.
- [96] C. L. Blake and C. J. Merz, "Uci repository of machine learning databases," 1998, 1998. [Online]. Available: <http://www.ics.uci.edu/~mlearn/MLRepository.html>
- [97] A. S. Georgiades, P. N. Belhumeur, and D. J. Kriegman, "From few to many: Generative models for recognition under variable pose and illumination," in *Proc. IEEE Int. Conf. Automatic Face and Gesture recognit. (cat. no. pr00580)*, 2000, pp. 277–284.
- [98] M. Ye and Y. Guo, "Deep triplet ranking networks for one-shot recognition," *arXiv preprint arXiv:1804.07275*, 2018.
- [99] A. Mehrotra and A. Dukkipati, "Generative adversarial residual pairwise networks for one shot learning," *arXiv preprint arXiv:1703.08033*, 2017.
- [100] L. Kaiser, O. Nachum, A. Roy, and S. Bengio, "Learning to remember rare events," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2017.
- [101] T. Scott, K. Ridgeway, and M. C. Mozer, "Adapted deep embeddings: A synthesis of methods for k-shot inductive transfer learning," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 76–85.
- [102] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu, "Learning fine-grained image similarity with deep ranking," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2014, pp. 1386–1393.
- [103] E. Ustinova and V. Lempitsky, "Learning deep embeddings with histogram loss," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2016, pp. 4170–4178.
- [104] Y. Bengio, S. Bengio, and J. Cloutier, *Learning a synaptic learning rule*. Université de Montréal, Département d'informatique et de recherche, 1990.
- [105] D. K. Naik and R. Mammone, "Meta-neural networks that learn by learning," in *Proc. Int. Joint Conf. Neural Net. (IJCNN)*, 1992.
- [106] S. Thrun and L. Pratt, *Learning to learn*. Springer Science & Business Media, 1998.
- [107] M. Andrychowicz, M. Denil, S. Gomez, M. W. Hoffman, D. Pfau, T. Schaul, B. Shillingford, and N. De Freitas, "Learning to learn by gradient descent by gradient descent," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2016, pp. 3981–3989.
- [108] K. Li and J. Malik, "Learning to optimize," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2017.
- [109] Y. Chen, M. W. Hoffman, S. G. Colmenarejo, M. Denil, T. P. Lillicrap, M. Botvinick, and N. de Freitas, "Learning to learn without gradient descent by gradient descent," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2017, pp. 748–756.
- [110] J. Goldberger, G. E. Hinton, S. T. Roweis, and R. R. Salakhutdinov, "Neighbourhood components analysis," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2005, pp. 513–520.
- [111] E. Triantafillou, R. Zemel, and R. Urtasun, "Few-shot learning through an information retrieval lens," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 2255–2265.
- [112] I. Tsochantidis, T. Joachims, T. Hofmann, and Y. Altun, "Large margin methods for structured and interdependent output variables," *J. Mach. Learn. Research*, vol. 6, no. Sep, pp. 1453–1484, 2005.
- [113] T. Hazan, J. Keshet, and D. A. McAllester, "Direct loss minimization for structured prediction," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2010, pp. 1594–1602.
- [114] B. Oreshkin, P. R. López, and A. Lacoste, "Tadam: Task dependent adaptive metric for improved few-shot learning," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 721–731.

- [115] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, "Film: Visual reasoning with a general conditioning layer," in *Proc. AAAI Conf. Artificial Intell. (AAAI)*, 2018, pp. 3942–3951.
- [116] C. Xing, N. Rostamzadeh, B. Oreshkin, and P. O. Pinheiro, "Adaptive cross-modal few-shot learning," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2019, pp. 4848–4858.
- [117] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Proc. Conf. Empirical Methods Nat. Language Process. (EMNLP)*, 2014, pp. 1532–1543.
- [118] F. Hao, J. Cheng, L. Wang, and J. Cao, "Instance-level embedding adaptation for few-shot learning," *IEEE Access*, vol. 7, pp. 100501–100511, 2019.
- [119] P. Hu, X. Sun, K. Saenko, and S. Sclaroff, "Weakly-supervised compositional feature aggregation for few-shot recognition," *arXiv preprint arXiv:1906.04833*, 2019.
- [120] X. Li, L. Yu, C.-W. Fu, M. Fang, and P.-A. Heng, "Revisiting metric learning for few-shot image classification," *Neurocomputing*, 2020.
- [121] Y. Zheng, R. Wang, J. Yang, L. Xue, and M. Hu, "Principal characteristic networks for few-shot learning," *J. Vis. Commun. Image Represent.*, vol. 59, pp. 563–573, 2019.
- [122] N. Dvornik, C. Schmid, and J. Mairal, "Diversity with cooperation: Ensemble methods for few-shot classification," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 3723–3731.
- [123] H. Prol, V. Dumoulin, and L. Herranz, "Cross-modulation networks for few-shot learning," *arXiv preprint arXiv:1812.00273*, 2018.
- [124] Q. Cai, Y. Pan, T. Yao, C. Yan, and T. Mei, "Memory matching networks for one-shot image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 4080–4088.
- [125] J. Weston, S. Chopra, and A. Bordes, "Memory networks," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2014.
- [126] L. Zhang, J. Liu, M. Luo, X. Chang, Q. Zheng, and A. G. Hauptmann, "Scheduled sampling for one-shot learning via matching network," *Pattern Recognit.*, vol. 96, p. 106962, 2019.
- [127] N. Hilliard, L. Phillips, S. Howland, A. Yankov, C. D. Corley, and N. O. Hodas, "Few-shot learning with metric-agnostic conditional embeddings," *arXiv preprint arXiv:1802.04376*, 2018.
- [128] X. Zhang, F. Sung, Y. Qiang, Y. Yang, and T. M. Hospedales, "Deep comparison: Relation columns for few-shot learning," *arXiv preprint arXiv:1811.07100*, 2018.
- [129] W. Li, J. Xu, J. Huo, L. Wang, Y. Gao, and J. Luo, "Distribution consistency based covariance metric networks for few-shot learning," in *Proc. AAAI Conf. Artificial Intell. (AAAI)*, vol. 33, 2019, pp. 8642–8649.
- [130] W. Li, L. Wang, J. Xu, J. Huo, Y. Gao, and J. Luo, "Revisiting local descriptor based image-to-class measure for few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 7260–7268.
- [131] H. Zhang and P. Koniusz, "Power normalizing second-order similarity network for few-shot learning," in *Proc. IEEE Winter Conf. Applica. Comput. Vis. (WACV)*, 2019, pp. 1185–1193.
- [132] B. Hui, P. Zhu, Q. Hu, and Q. Wang, "Self-attention relation network for few-shot learning," in *Proc. IEEE Int. Conf. Multimedia Expo Workshops (ICMEW)*, 2019, pp. 198–203.
- [133] V. Garcia and J. Bruna, "Few-shot learning with graph neural networks," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2018.
- [134] J. Kim, T. Kim, S. Kim, and C. D. Yoo, "Edge-labeling graph neural network for few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 11–20.
- [135] Y. Wang, X.-M. Wu, Q. Li, J. Gu, W. Xiang, L. Zhang, and V. O. Li, "Large margin few-shot learning," *arXiv preprint arXiv:1807.02872*, 2018.
- [136] H. Li, D. Eigen, S. Dodge, M. Zeiler, and X. Wang, "Finding task-relevant features for few-shot learning by category traversal," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 1–10.
- [137] Z. Li, F. Zhou, F. Chen, and H. Li, "Meta-sgd: Learning to learn quickly for few-shot learning," *arXiv preprint arXiv:1707.09835*, 2017.
- [138] Y. Lee and S. Choi, "Gradient-based meta-learning with learned layerwise metric and subspace," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 2927–2936.
- [139] A. Nichol and J. Schulman, "Reptile: a scalable metalearning algorithm," *arXiv preprint arXiv:1803.02999*, vol. 2, p. 2, 2018.
- [140] F. Zhou, B. Wu, and Z. Li, "Deep meta-learning: Learning to learn in the concept space," *arXiv preprint arXiv:1802.03596*, 2018.
- [141] E. Grant, C. Finn, S. Levine, T. Darrell, and T. Griffiths, "Recasting gradient-based meta-learning as hierarchical bayes," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2018.
- [142] C. Finn, K. Xu, and S. Levine, "Probabilistic model-agnostic meta-learning," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 9516–9527.
- [143] J. Yoon, T. Kim, O. Dia, S. Kim, Y. Bengio, and S. Ahn, "Bayesian model-agnostic meta-learning," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 7332–7342.
- [144] M. A. Jamal and G.-J. Qi, "Task agnostic meta-learning for few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 11719–11727.
- [145] A. A. Rusu, D. Rao, J. Sygnowski, O. Vinyals, R. Pascanu, S. Osindero, and R. Hadsell, "Meta-learning with latent embedding optimization," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.
- [146] X. Jiang, M. Havvaei, F. Varno, G. Chartrand, N. Chapados, and S. Matwin, "Learning to learn with conditional class dependencies," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.
- [147] Q. Liu and D. Wang, "Stein variational gradient descent: A general purpose bayesian inference algorithm," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2016, pp. 2378–2386.
- [148] J. Nie, N. Xu, M. Zhou, G. Yan, and Z. Wei, "3d model classification based on few-shot learning," *Neurocomputing*, vol. 398, pp. 539–546, 2020.
- [149] Q. Sun, Y. Liu, T.-S. Chua, and B. Schiele, "Meta-transfer learning for few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 403–412.
- [150] L. Bertinetto, J. F. Henriques, J. Valmadre, P. Torr, and A. Vedaldi, "Learning feed-forward one-shot learners," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2016, pp. 523–531.
- [151] Y.-X. Wang and M. Hebert, "Learning to learn: Model regression networks for easy small sample learning," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*. Springer, 2016, pp. 616–634.
- [152] S. Gidaris and N. Komodakis, "Dynamic few-shot visual learning without forgetting," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 4367–4375.
- [153] S. Qiao, C. Liu, W. Shen, and A. L. Yuille, "Few-shot image recognition by predicting parameters from activations," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 7229–7238.
- [154] H. Qi, M. Brown, and D. G. Lowe, "Low-shot learning with imprinted weights," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 5822–5830.
- [155] F. Zhao, J. Zhao, S. Yan, and J. Feng, "Dynamic conditional networks for few-shot learning," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 19–35.
- [156] T. Wu, J. Peurifoy, I. L. Chuang, and M. Tegmark, "Meta-learning autoencoders for few-shot prediction," *arXiv preprint arXiv:1807.09912*, 2018.
- [157] S. Gidaris and N. Komodakis, "Generating classification weights with gnn denoising autoencoders for few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 21–30.
- [158] J. Gordon, J. Bronskill, M. Bauer, S. Nowozin, and R. E. Turner, "Meta-learning probabilistic inference for prediction," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.
- [159] L. Bertinetto, J. F. Henriques, P. H. Torr, and A. Vedaldi, "Meta-learning with differentiable closed-form solvers," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.
- [160] K. Lee, S. Maji, A. Ravichandran, and S. Soatto, "Meta-learning with differentiable convex optimization," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 10657–10665.
- [161] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2014.
- [162] T. Munkhdalai, X. Yuan, S. Mehri, and A. Trischler, "Rapid adaptation with conditionally shifted neurons," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 3664–3673.
- [163] T. Munkhdalai and A. Trischler, "Metalearning with hebbian fast weights," *arXiv preprint arXiv:1807.05076*, 2018.
- [164] H.-J. Ye, H. Hu, D.-C. Zhan, and F. Sha, "Few-shot learning via embedding adaptation with set-to-set functions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020.
- [165] D. O. Hebb, *The organization of behavior: a neuropsychological theory*. J. Wiley; Chapman & Hall, 1949.
- [166] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. R. Salakhutdinov, and A. J. Smola, "Deep sets," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 3391–3401.

- [167] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2017.
- [168] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 5998–6008.
- [169] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap, "Meta-learning with memory-augmented neural networks," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2016, pp. 1842–1850.
- [170] P. Shyam, S. Gupta, and A. Dukkipati, "Attentive recurrent comparators," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2017, pp. 3173–3181.
- [171] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel, "A simple neural attentive meta-learner," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2018.
- [172] T. Ramalho and M. Garnelo, "Adaptive posterior learning: few-shot learning with a surprise-based memory module," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.
- [173] A. Graves, G. Wayne, and I. Danihelka, "Neural turing machines," *arXiv preprint arXiv:1410.5401*, 2014.
- [174] S. Hochreiter, Y. Bengio, P. Frasconi, J. Schmidhuber *et al.*, "Gradient flow in recurrent nets: the difficulty of learning long-term dependencies." A field guide to dynamical recurrent neural networks. IEEE Press, 2001.
- [175] Y. Zhang and Q. Yang, "A survey on multi-task learning," *arXiv preprint arXiv:1707.08114*, 2017.
- [176] W. Yan, J. Yap, and G. Mori, "Multi-task transfer methods to improve one-shot learning for multimedia event detection," in *Proc. British Mach. Vis. Conf. (BMVC)*, 2015, pp. 3701–3713.
- [177] R. Zhang, T. Che, Z. Ghahramani, Y. Bengio, and Y. Song, "Metagan: An adversarial approach to few-shot learning," in *Proc. Advances Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 2365–2374.
- [178] Z. Hu, X. Li, C. Tu, Z. Liu, and M. Sun, "Few-shot charge prediction with discriminative legal attributes," in *Proc. Int. Conf. Computational Linguistics*, 2018, pp. 487–498.
- [179] J.-C. Su, S. Maji, and B. Hariharan, "Boosting supervision with self-supervision for few-shot learning," *arXiv preprint arXiv:1906.07079*, 2019.
- [180] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2015, pp. 1422–1430.
- [181] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*. Springer, 2016, pp. 69–84.
- [182] H.-Y. Tung, H.-W. Tung, E. Yumer, and K. Fragkiadaki, "Self-supervised learning of motion capture," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 5236–5246.
- [183] B. Fernando, H. Bilen, E. Gavves, and S. Gould, "Self-supervised video representation learning with odd-one-out networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 3636–3645.
- [184] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2018.
- [185] M. Noroozi, A. Vinjimoor, P. Favaro, and H. Pirsiavash, "Boosting self-supervised learning via knowledge transfer," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 9359–9367.
- [186] T. Nathan Mundhenk, D. Ho, and B. Y. Chen, "Improvements to context based self-supervised learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 9339–9348.
- [187] D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song, "Using self-supervised learning can improve model robustness and uncertainty," in *Proc. Advances Neural Inf. Process. Syst. (NeurIPS)*, 2019, pp. 15 663–15 674.
- [188] A. Kolesnikov, X. Zhai, and L. Beyer, "Revisiting self-supervised visual representation learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 1920–1929.
- [189] L. Jing and Y. Tian, "Self-supervised visual feature learning with deep neural networks: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.
- [190] H. Wang, X. Wang, W. Xiong, M. Yu, X. Guo, S. Chang, and W. Y. Wang, "Self-supervised learning for contextualized extractive summarization," in *Proc. Annu. Meet. Assoc. Computational Linguistics (ACL)*, 2019, pp. 2221–2227.
- [191] D. Kang, A. Balakrishnan, P. Shah, P. A. Crook, Y.-L. Boureau, and J. Weston, "Recommendation as a communication game: Self-supervised bot-play for goal-oriented dialogue," in *Proc. Conf. Empirical Methods Nat. Language Process. (EMNLP)*, 2019, pp. 1951–1961.
- [192] X. Wang, Q. Huang, A. Celikyilmaz, J. Gao, D. Shen, Y.-F. Wang, W. Y. Wang, and L. Zhang, "Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 6629–6638.
- [193] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "Albert: A lite bert for self-supervised learning of language representations," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2020.
- [194] A. Dosovitskiy, J. T. Springenberg, M. Riedmiller, and T. Brox, "Discriminative unsupervised feature learning with convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2014, pp. 766–774.
- [195] S. Gidaris, A. Bursuc, N. Komodakis, P. Pérez, and M. Cord, "Boosting few-shot visual learning with self-supervision," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 8059–8068.
- [196] X. Li, Q. Sun, Y. Liu, Q. Zhou, S. Zheng, T.-S. Chua, and B. Schiele, "Learning to self-train for semi-supervised few-shot classification," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2019, pp. 10 276–10 286.
- [197] P. Mangla, N. Kumari, A. Sinha, M. Singh, B. Krishnamurthy, and V. N. Balasubramanian, "Charting the right manifold: Manifold mixup for few-shot learning," in *Proc. IEEE Winter Conf. Application Comput. Vis. (WACV)*, 2020, pp. 2218–2227.
- [198] V. N. Vapnik, "An overview of statistical learning theory," *IEEE Trans. Neural Net.*, vol. 10, no. 5, pp. 988–999, 1999.
- [199] Y. Liu, J. Lee, M. Park, S. Kim, E. Yang, S. J. Hwang, and Y. Yang, "Learning to propagate labels: Transductive propagation network for few-shot learning," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.
- [200] Z. Xu, L. Zhu, and Y. Yang, "Few-shot object recognition from machine-labeled web images," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 1164–1172.
- [201] M. Bauer, M. Rojas-Carulla, J. B. Świątkowski, B. Schölkopf, and R. E. Turner, "Discriminative k-shot learning using probabilistic models," *arXiv preprint arXiv:1706.00326*, 2017.
- [202] L. Franceschi, P. Frasconi, S. Salzo, R. Grazzi, and M. Pontil, "Bilevel programming for hyperparameter optimization and meta-learning," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 1568–1577.
- [203] A. Kimura, Z. Ghahramani, K. Takeuchi, T. Iwata, and N. Ueda, "Few-shot learning of neural networks from scratch by pseudo example optimization," *arXiv preprint arXiv:1802.03039*, 2018.
- [204] Y. Lifchitz, Y. Avrithis, S. Picard, and A. Bursuc, "Dense classification and implanting for few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 9258–9267.
- [205] H. Zhang, J. Zhang, and P. Koniusz, "Few-shot learning via saliency-guided hallucination of samples," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 2770–2779.
- [206] W.-Y. Chen, Y.-C. Liu, Z. Kira, Y.-C. F. Wang, and J.-B. Huang, "A closer look at few-shot classification," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.
- [207] Y. Wang, W.-L. Chao, K. Q. Weinberger, and L. van der Maaten, "Simpleshot: Revisiting nearest-neighbor classification for few-shot learning," *arXiv preprint arXiv:1911.04623*, 2019.
- [208] R. Boney and A. Ilin, "Semi-supervised and active few-shot learning with prototypical networks," *arXiv preprint arXiv:1711.10856*, 2017.
- [209] M. Ren, E. Triantafillou, S. Ravi, J. Snell, K. Swersky, J. B. Tenenbaum, H. Larochelle, and R. S. Zemel, "Meta-learning for semi-supervised few-shot classification," *arXiv preprint arXiv:1803.00676*, 2018.
- [210] K. R. Allen, E. Shelhamer, H. Shin, and J. B. Tenenbaum, "Infinite mixture prototypes for few-shot learning," *arXiv preprint arXiv:1902.04552*, 2019.
- [211] J. Wang, K. Chen, L. Shou, S. Wu, and S. Mehrotra, "Semi-supervised few-shot learning for dual question-answer extraction," *arXiv preprint arXiv:1904.03898*, 2019.
- [212] T. Ma and A. Zhang, "Affinitynet: semi-supervised few-shot learning for disease type prediction," in *Proc. AAAI Conf. Artificial Intell. (AAAI)*, vol. 33, 2019, pp. 1069–1076.
- [213] K. Hsu, S. Levine, and C. Finn, "Unsupervised learning via meta-learning," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.

- [214] D. Berthelot, C. Raffel, A. Roy, and I. Goodfellow, "Understanding and improving interpolation in autoencoders via an adversarial regularizer," *arXiv preprint arXiv:1807.07543*, 2018.
- [215] S. Khodadadeh, L. Boloni, and M. Shah, "Unsupervised meta-learning for few-shot image classification," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2019, pp. 10 132–10 142.
- [216] A. Antoniou and A. Storkey, "Assume, augment and learn: Unsupervised few-shot meta-learning via random labels and data augmentation," *arXiv preprint arXiv:1902.09884*, 2019.
- [217] V. M. Patel, R. Gopalan, R. Li, and R. Chellappa, "Visual domain adaptation: A survey of recent advances," *IEEE Signal Process. Mag.*, vol. 32, no. 3, pp. 53–69, 2015.
- [218] S. Motian, Q. Jones, S. Iranmanesh, and G. Doretto, "Few-shot adversarial domain adaptation," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 6670–6680.
- [219] Z. Luo, Y. Zou, J. Hoffman, and L. F. Fei-Fei, "Label efficient learning of transferable representations across domains and tasks," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 165–177.
- [220] N. Dong and E. P. Xing, "Domain adaption in one-shot learning," in *Proc. Joint Eur. Conf. Mach. Learn. Know. Dis. in Databases*. Springer, 2018, pp. 573–588.
- [221] B. Kang and J. Feng, "Transferable meta learning across domains," in *Proc. Conf. Uncertainty in Artificial Intell.*, 2018, pp. 177–187.
- [222] D. Sahoo, H. Le, C. Liu, and S. C. Hoi, "Meta-learning with domain adaptation for few-shot learning under domain shift," *submitted to ICLR 2019*.
- [223] H.-Y. Tseng, H.-Y. Lee, J.-B. Huang, and M.-H. Yang, "Cross-domain few-shot classification via learned feature-wise transformation," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2020.
- [224] J. Lu, Z. Cao, K. Wu, G. Zhang, and C. Zhang, "Boosting few-shot image recognition via domain alignment prototypical networks," in *Proc. IEEE Int. Conf. Tools Artificial Intell. (ICTAI)*, 2018, pp. 260–264.
- [225] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS) Workshop on Deep Learning and Unsupervised Feature Learning*, 2011, pp. 165–177.
- [226] K. Soomro, A. R. Zamir, and M. Shah, "Ucf101: A dataset of 101 human actions classes from videos in the wild," *arXiv preprint arXiv:1212.0402*, 2012.
- [227] Y. Le Cun, L. D. Jackel, B. Boser, J. S. Denker, H. P. Graf, I. Guyon, D. Henderson, R. E. Howard, and W. Hubbard, "Handwritten digit recognition: Applications of neural network chips and automatic learning," *IEEE Commun. Mag.*, vol. 27, no. 11, pp. 41–46, 1989.
- [228] K. Saenko, B. Kulis, M. Fritz, and T. Darrell, "Adapting visual category models to new domains," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*. Springer, 2010, pp. 213–226.
- [229] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan, "Deep hashing network for unsupervised domain adaptation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 5018–5027.
- [230] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, "3d object representations for fine-grained categorization," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV) Workshop*, 2013, pp. 554–561.
- [231] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 6, pp. 1452–1464, 2017.
- [232] G. Van Horn, O. Mac Aodha, Y. Song, Y. Cui, C. Sun, A. Shepard, H. Adam, P. Perona, and S. Belongie, "The inaturalist species classification and detection dataset," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 8769–8778.
- [233] L. Zhu and Y. Yang, "Compound memory networks for few-shot video classification," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 751–766.
- [234] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio, "An empirical investigation of catastrophic forgetting in gradient-based neural networks," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2014.
- [235] X. Shi, L. Salewski, M. Schiegg, Z. Akata, and M. Welling, "Relational generalized few-shot learning," *arXiv preprint arXiv:1907.09557*, 2019.
- [236] M. Ren, R. Liao, E. Fetaya, and R. Zemel, "Incremental few-shot learning with attention attractor networks," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2019, pp. 5276–5286.
- [237] S. Rahman, S. Khan, and F. Porikli, "A unified approach for conventional zero-shot, generalized zero-shot, and few-shot learning," *IEEE Trans. Image Process.*, vol. 27, no. 11, pp. 5652–5667, 2018.
- [238] E. Schonfeld, S. Ebrahimi, S. Sinha, T. Darrell, and Z. Akata, "Generalized zero-and few-shot learning via aligned variational autoencoders," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 8247–8255.
- [239] Y. Huang and L. Wang, "Acmm: Aligned cross-modal memory for few-shot image and sentence matching," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 5774–5783.
- [240] Y. Huang, Y. Long, and L. Wang, "Few-shot image and sentence matching via gated visual-semantic embedding," in *Proc. AAAI Conf. Artificial Intell. (AAAI)*, vol. 33, 2019, pp. 8489–8496.
- [241] R. Elloff, H. A. Engelbrecht, and H. Kamper, "Multimodal one-shot learning of speech and images," in *Proc. IEEE Int. Conf. Acoustics Speech Sig. Process. (ICASSP)*, 2019, pp. 8623–8627.
- [242] F. Pahde, M. Nabi, T. Klein, and P. Jähnichen, "Discriminative hallucination for multi-modal few-shot learning," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, 2018, pp. 156–160.
- [243] F. Pahde, P. Jähnichen, T. Klein, and M. Nabi, "Cross-modal hallucination for few-shot fine-grained recognition," *arXiv preprint arXiv:1806.05147*, 2018.
- [244] Y.-H. H. Tsai and R. Salakhutdinov, "Improving one-shot learning through fusing side information," *arXiv preprint arXiv:1710.08347*, 2017.
- [245] M. P. Fortin and B. Chaib-draa, "Few-shot learning with contextual cueing for object recognition in complex scenes," *arXiv preprint arXiv:1912.06679*, 2019.
- [246] R. Vuorio, S.-H. Sun, H. Hu, and J. J. Lim, "Multimodal model-agnostic meta-learning via task-aware modulation," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2019, pp. 1–12.
- [247] F. Pahde, O. Ostapenko, P. J. Hnichen, T. Klein, and M. Nabi, "Self-paced adversarial training for multimodal few-shot learning," in *Proc. IEEE Winter Conf. Applica. Comput. Vis. (WACV)*, 2019, pp. 218–226.
- [248] A. Li, T. Luo, Z. Lu, T. Xiang, and L. Wang, "Large-scale few-shot learning: Knowledge transfer with class hierarchy," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 7212–7220.
- [249] P. Tokmakov, Y.-X. Wang, and M. Hebert, "Learning compositional representations for few-shot recognition," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 6372–6381.
- [250] S. Yan, S. Zhang, X. He *et al.*, "A dual attention network with semantic embedding for few-shot learning," in *Proc. AAAI Conf. Artificial Intell. (AAAI)*, vol. 33, 2019, pp. 9079–9086.
- [251] A. Alfassy, L. Karlinsky, A. Aides, J. Shtok, S. Harary, R. Feris, R. Giryes, and A. M. Bronstein, "Laso: Label-set operations networks for multi-label few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 6548–6557.
- [252] A. Ali-Gombe, E. Elyan, Y. Savoye, and C. Jayne, "Few-shot classifier gan," in *Proc. Int. Joint Conf. Neural Net. (IJCNN)*, 2018, pp. 1–8.
- [253] F. Pahde, M. Puscas, J. Wolff, T. Klein, N. Sebe, and M. Nabi, "Low-shot learning from imaginary 3d model," in *Proc. IEEE Winter Conf. Applica. Comput. Vis. (WACV)*, 2019, pp. 978–985.
- [254] X.-S. Wei, P. Wang, L. Liu, C. Shen, and J. Wu, "Piecewise classifier mappings: Learning fine-grained learners for novel categories with few examples," *IEEE Trans. Image Process.*, vol. 28, no. 12, pp. 6116–6125, 2019.
- [255] H. Huang, J. Zhang, J. Zhang, J. Xu, and Q. Wu, "Low-rank pairwise alignment bilinear network for few-shot fine-grained image classification," *arXiv preprint arXiv:1908.01313*, 2019.
- [256] B. Liu, X. Yu, A. Yu, P. Zhang, G. Wan, and R. Wang, "Deep few-shot learning for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sensing*, vol. 57, no. 4, pp. 2290–2304, 2018.
- [257] C. Zhang, J. Yue, and Q. Qin, "Deep quadruplet network for hyperspectral image classification with a small number of samples," *Remote Sensing*, vol. 12, no. 4, p. 647, 2020.
- [258] A. Oliver, X. Lladó, A. Torrent, and J. Martí, "One-shot segmentation of breast, pectoral muscle, and background in digitised mammograms," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, 2014, pp. 912–916.
- [259] A. Shaban, S. Bansal, Z. Liu, I. Essa, and B. Boots, "One-shot learning for semantic segmentation," *arXiv preprint arXiv:1709.03410*, 2017.
- [260] K. Rakelly, E. Shelhamer, T. Darrell, A. A. Efros, and S. Levine, "Few-shot segmentation propagation with guided networks," *arXiv preprint arXiv:1806.07373*, 2018.

- [261] N. Dong and E. Xing, "Few-shot semantic segmentation with prototype learning," in *Proc. British Mach. Vis. Conf. (BMVC)*, vol. 3, no. 4, 2018.
- [262] X. Zhang, Y. Wei, Y. Yang, and T. Huang, "Sg-one: Similarity guidance network for one-shot semantic segmentation," *arXiv preprint arXiv:1810.09091*, 2018.
- [263] Z. Dong, R. Zhang, X. Shao, and H. Zhou, "Multi-scale discriminative location-aware network for few-shot semantic segmentation," in *Proc. IEEE Annu. Comput. Soft. App. Conf. (COMPSAC)*, vol. 2, 2019, pp. 42–47.
- [264] M. Siam, B. N. Oreshkin, and M. Jagersand, "Amp: Adaptive masked proxies for few-shot segmentation," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 5249–5258.
- [265] T. Hu, P. Yang, C. Zhang, G. Yu, Y. Mu, and C. G. Snoek, "Attention-based multi-context guiding for few-shot semantic segmentation," in *Proc. AAAI Conf. Artificial Intell. (AAAI)*, vol. 33, 2019, pp. 8441–8448.
- [266] C. Zhang, G. Lin, F. Liu, R. Yao, and C. Shen, "Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 5217–5226.
- [267] C. Michaelis, M. Bethge, and A. Ecker, "One-shot segmentation in clutter," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 3549–3558.
- [268] A. K. Bhunia, A. K. Bhunia, S. Ghose, A. Das, P. P. Roy, and U. Pal, "A deep one-shot network for query-based logo retrieval," *Pattern Recognit.*, vol. 96, pp. 1–10, 2019.
- [269] I. Ustyuzhaninov, C. Michaelis, W. Brendel, and M. Bethge, "One-shot texture segmentation," *arXiv preprint arXiv:1807.02654*, 2018.
- [270] K. Zhu, W. Zhai, Z.-J. Zha, and Y. Cao, "One-shot texture retrieval with global context metric," in *Proc. Int. Joint Conf. Artificial Intell. (IJCAI)*, 2019, pp. 4461–4467.
- [271] A. K. Mondal, J. Dolz, and C. Desrosiers, "Few-shot 3d multi-modal medical image segmentation using generative adversarial learning," *arXiv preprint arXiv:1810.12241*, 2018.
- [272] A. Zhao, G. Balakrishnan, F. Durand, J. V. Guttag, and A. V. Dalca, "Data augmentation using learned transforms for one-shot medical image segmentation," *arXiv preprint arXiv:1902.09383*, 2019.
- [273] J. Dietlmeier, K. McGuinness, S. Rugonyi, T. Wilson, A. Nuttall, and N. E. O'Connor, "Few-shot hypercolumn-based mitochondria segmentation in cardiac and outer hair cells in focused ion beam-scanning electron microscopy (fib-sem) data," *Pattern Recognit. Letters*, vol. 128, pp. 521–528, 2019.
- [274] A. G. Roy, S. Siddiqui, S. Pölsterl, N. Navab, and C. Wachinger, "'squeeze & excite' guided few-shot segmentation of volumetric images," *Medical Image Anal.*, vol. 59, pp. 1–12, 2020.
- [275] X. Dong, L. Zheng, F. Ma, Y. Yang, and D. Meng, "Few-example object detection with model communication," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 7, pp. 1641–1654, 2018.
- [276] J. Chen, Y. Liu, Y. Liu, S. Wang, and S. Chen, "A few-shot learning framework for air vehicle detection by similarity embedding," in *Proc. Int. Conf. Graphics Image Process. (ICGIP)*, 2019.
- [277] Q. Fan, W. Zhuo, and Y.-W. Tai, "Few-shot object detection with attention-rpn and multi-relation detector," *arXiv preprint arXiv:1908.01998*, 2019.
- [278] L. Karlinsky, J. Shtok, S. Harary, E. Schwartz, A. Aides, R. Feris, R. Giryes, and A. M. Bronstein, "Repmet: Representative-based metric learning for classification and few-shot object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 5197–5206.
- [279] B. Kang, Z. Liu, X. Wang, F. Yu, J. Feng, and T. Darrell, "Few-shot object detection via feature reweighting," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 8420–8429.
- [280] J. Sun, Y. Zhu, and S. Jiang, "One-shot learning for rgb-d handheld object recognition," in *Proc. Int. Conf. Internet Multimedia Comput. Service*, 2018, pp. 1–6.
- [281] S. Benaim and L. Wolf, "One-shot unsupervised cross domain translation," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 2104–2114.
- [282] L. Clouâtre and M. Demers, "Figr: Few-shot image generation with reptile," *arXiv preprint arXiv:1901.02199*, 2019.
- [283] M.-Y. Liu, X. Huang, A. Mallya, T. Karras, T. Aila, J. Lehtinen, and J. Kautz, "Few-shot unsupervised image-to-image translation," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 10551–10560.
- [284] Y. Gao, Y. Guo, Z. Lian, Y. Tang, and J. Xiao, "Artistic glyph image synthesis via one-stage few-shot learning," *ACM Trans. Graphics*, vol. 38, no. 6, pp. 1–12, 2019.
- [285] Y. Hong, J. Zhang, L. Niu, and L. Zhang, "Matchinggan: Matching-based few-shot image generation," *arXiv preprint arXiv:2003.03497*, 2020.
- [286] Y.-X. Wang, L. Gui, and M. Hebert, "Few-shot hash learning for image retrieval," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV) Workshop*, 2017, pp. 1228–1237.
- [287] Y. Yu, G. Liu, and J.-M. Odobez, "Improving few-shot user-specific gaze adaptation via gaze redirection synthesis," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 11937–11946.
- [288] S. Li, J. Shi, W. Song, A. Hao, and H. Qin, "Few-shot learning for monocular depth estimation based on local object relationship," in *Proc. IEEE Int. Conf. Tools Artificial Intell. (ICTAI)*, 2019, pp. 1221–1228.
- [289] D. Wertheimer and B. Hariharan, "Few-shot learning with localization in realistic settings," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 6558–6567.
- [290] A. Dornadula, A. Narcomey, R. Krishna, M. Bernstein, and F.-F. Li, "Visual relationships as functions: Enabling few-shot scene graph prediction," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV) Workshop*, 2019.
- [291] L. Xiang, X. Jin, G. Ding, J. Han, and L. Li, "Incremental few-shot learning for pedestrian attribute recognition," *arXiv preprint arXiv:1906.00330*, 2019.
- [292] T. Xu, J. Li, H. Wu, H. Yang, X. Gu, and Y. Chen, "Feature space regularization for person re-identification with one sample," in *Proc. IEEE Int. Conf. Tools Artificial Intell. (ICTAI)*, 2019, pp. 1463–1470.
- [293] S. Yoo, H. Bahng, S. Chung, J. Lee, J. Chang, and J. Choo, "Coloring with limited data: Few-shot colorization via memory augmented networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 11283–11292.
- [294] S. McDonagh, S. Parisot, F. Zhou, X. Zhang, A. Leonardis, Z. Li, and G. Slabaugh, "Formulating camera-adaptive color constancy as a few-shot meta-learning problem," *arXiv*, pp. arXiv–1811, 2018.
- [295] S. Khodadadeh, L. Böloni, and M. Shah, "Unsupervised meta-learning for few-shot image and video classification," *arXiv preprint arXiv:1811.11819*, 2018.
- [296] K. Cao, J. Ji, Z. Cao, C.-Y. Chang, and J. C. Niebles, "Few-shot video classification via temporal alignment," *arXiv preprint arXiv:1906.11415*, 2019.
- [297] X. Li, S. Qin, K. Xu, and Z. Hu, "One-shot learning gesture recognition based on evolution of discrimination with successive memory," in *Proc. IEEE Int. Conf. Intell. Robotic Control Eng. (IRCE)*, 2018, pp. 263–269.
- [298] Z. Lu, S. Qin, X. Li, L. Li, and D. Zhang, "One-shot learning hand gesture recognition based on modified 3d convolutional neural networks," *Mach. Vis. App.*, vol. 30, no. 7–8, pp. 1157–1180, 2019.
- [299] A. Mishra, V. K. Verma, M. S. K. Reddy, S. Arulkumar, P. Rai, and A. Mittal, "A generative approach to zero-shot and few-shot action recognition," in *Proc. IEEE Winter Conf. Applica. Comput. Vis. (WACV)*, 2018, pp. 372–380.
- [300] B. Xu, H. Ye, Y. Zheng, H. Wang, T. Luwang, and Y.-G. Jiang, "Dense dilated network for few shot action recognition," in *Proc. ACM Int. Conf. Multimedia Retrieval*, 2018, pp. 379–387.
- [301] M. Bishay, G. Zoumpourlis, and I. Patras, "Tarn: Temporal attentive relation network for few-shot and zero-shot action recognition," *arXiv preprint arXiv:1907.09021*, 2019.
- [302] S. K. Dwivedi, V. Gupta, R. Mitra, S. Ahmed, and A. Jain, "Protogan: Towards few shot learning for action recognition," *arXiv preprint arXiv:1909.07945*, 2019.
- [303] H. Coskun, Z. Zia, B. Tekin, F. Bogo, N. Navab, F. Tombari, and H. Sawhney, "Domain-specific priors and meta learning for low-shot first-person action recognition," *arXiv preprint arXiv:1907.09382*, 2019.
- [304] R. Memmesheimer, N. Theisen, and D. Paulus, "Signal level deep metric learning for multimodal one-shot action recognition," *arXiv preprint arXiv:2004.11085*, 2020.
- [305] H. Yang, X. He, and F. Porikli, "One-shot action localization by learning sequence matching network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 1450–1459.
- [306] W. Goo and S. Nieku, "One-shot learning of multi-step tasks from observation via activity localization in auxiliary video," in *Proc. Int. Conf. Robotics Auto. (ICRA)*, 2019, pp. 7755–7761.
- [307] H. Xu, B. Kang, X. Sun, J. Feng, K. Saenko, and T. Darrell, "Similarity r-c3d for few-shot temporal activity detection," *arXiv preprint arXiv:1812.10000*, 2018.

- [308] L.-Y. Gui, Y.-X. Wang, D. Ramanan, and J. M. Moura, "Few-shot human motion prediction via meta-learning," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 432–450.
- [309] S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool, "One-shot video object segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 221–230.
- [310] H. Xiao, B. Kang, Y. Liu, M. Zhang, and J. Feng, "Online meta adaptation for fast video object segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019.
- [311] N. Inoue and K. Shinoda, "Few-shot adaptation for multimedia semantic indexing," in *Proc. ACM Int. Conf. Multimedia (AMM)*, 2018, pp. 1110–1118.
- [312] J. Lee, D. Ramanan, and R. Girdhar, "Metapix: Few-shot video retargeting," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2020.
- [313] E. Zakharov, A. Shysheya, E. Burkov, and V. Lempitsky, "Few-shot adversarial learning of realistic neural talking head models," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 9459–9468.
- [314] L. Wu, Y. Wang, H. Yin, M. Wang, and L. Shao, "Few-shot deep adversarial learning for video-based person re-identification," *IEEE Trans. Image Process.*, vol. 29, pp. 1233–1245, 2019.
- [315] E. Park and A. C. Berg, "Meta-tracker: Fast and robust online adaptation for visual object trackers," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 569–585.
- [316] I. Mason, S. Starke, H. Zhang, H. Bilen, and T. Komura, "Few-shot learning of homogeneous human locomotion styles," in *Comput. Graphics Forum*, vol. 37, no. 7, 2018, pp. 143–153.
- [317] L. Yan, Y. Zheng, and J. Cao, "Few-shot learning for short text classification," *Multimedia Tools App.*, vol. 77, no. 22, pp. 29 799–29 810, 2018.
- [318] M. Yu, X. Guo, J. Yi, S. Chang, S. Potdar, Y. Cheng, G. Tesauro, H. Wang, and B. Zhou, "Diverse few-shot text classification with multiple metrics," in *Proc. Conf. North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2018, pp. 1206–1215.
- [319] W. Liu, X. Chang, Y. Yan, Y. Yang, and A. G. Hauptmann, "Few-shot text and image classification via analogical transfer learning," *ACM Trans. Intell. Syst. Techn. (TIST)*, vol. 9, no. 6, pp. 1–20, 2018.
- [320] X. Jiang, M. Havaei, G. Chartrand, H. Chouaib, T. Vincent, A. Jesson, N. Chapados, and S. Matwin, "On the importance of attention in meta-learning for few-shot text classification," *arXiv preprint arXiv:1806.00852*, 2018.
- [321] A. Rios and R. Kavuluru, "Few-shot and zero-shot multi-label learning for structured label spaces," in *Proc. Conf. Empirical Methods Nat. Language Process. (EMNLP)*, vol. 2018, 2018, p. 3132.
- [322] K. Bailey and S. Chopra, "Few-shot text classification with pre-trained word embeddings and a human in the loop," *arXiv preprint arXiv:1804.02063*, 2018.
- [323] Y. Bao, M. Wu, S. Chang, and R. Barzilay, "Few-shot text classification with distributional signatures," *arXiv preprint arXiv:1908.06039*, 2019.
- [324] N. Zhang, Z. Sun, S. Deng, J. Chen, and H. Chen, "Improving few-shot text classification via pretrained language representations," *arXiv preprint arXiv:1908.08788*, 2019.
- [325] R. Geng, B. Li, Y. Li, X. Zhu, P. Jian, and J. Sun, "Induction networks for few-shot text classification," in *Proc. Conf. Empirical Methods Nat. Language Process. (EMNLP)*, 2019, pp. 3895–3904.
- [326] S. Sun, Q. Sun, K. Zhou, and T. Lv, "Hierarchical attention prototypical networks for few-shot text classification," in *Proc. Conf. Empirical Methods Nat. Language Process. (EMNLP)*, 2019, pp. 476–485.
- [327] C. Pan, J. Huang, J. Gong, and X. Yuan, "Few-shot transfer learning for text classification with lightweight word embedding based models," *IEEE Access*, vol. 7, pp. 53 296–53 304, 2019.
- [328] T. Schick and H. Schütze, "Exploiting cloze questions for few-shot text classification and natural language inference," *arXiv preprint arXiv:2001.07676*, 2020.
- [329] V. Vlasov, A. Drissner-Schmid, and A. Nichol, "Few-shot generalization across dialogue tasks," *arXiv preprint arXiv:1811.11707*, 2018.
- [330] A. Madotto, Z. Lin, C.-S. Wu, and P. Fung, "Personalizing dialogue agents via meta-learning," in *Proc. Annu. Meet. Assoc. Computational Linguistics (ACL)*, 2019, pp. 5454–5459.
- [331] K. Qian and Z. Yu, "Domain adaptive dialog generation via meta learning," in *Proc. Annu. Meet. Assoc. Computational Linguistics (ACL)*, 2019, pp. 2639–2649.
- [332] X. Han, H. Zhu, P. Yu, Z. Wang, Y. Yao, Z. Liu, and M. Sun, "Fewrel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation," in *Proc. Conf. Empirical Methods Nat. Language Process. (EMNLP)*, 2018, pp. 4803–4809.
- [333] W. Xiong, M. Yu, S. Chang, X. Guo, and W. Y. Wang, "One-shot relational learning for knowledge graphs," in *Proc. Conf. Empirical Methods Nat. Language Process. (EMNLP)*, 2018, pp. 1980–1990.
- [334] C. Zhang, H. Yao, C. Huang, M. Jiang, Z. Li, and N. V. Chawla, "Few-shot knowledge graph completion," *arXiv preprint arXiv:1911.11298*, 2019.
- [335] Z.-X. Ye and Z.-H. Ling, "Multi-level matching and aggregation network for few-shot relation classification," in *Proc. Annu. Meet. Assoc. Computational Linguistics (ACL)*, 2019, pp. 2872–2881.
- [336] M. Chen, W. Zhang, W. Zhang, Q. Chen, and H. Chen, "Meta relational learning for few-shot link prediction in knowledge graphs," in *Proc. Conf. Empirical Methods Nat. Language Process. (EMNLP)*, 2019, pp. 4208–4217.
- [337] A. K. Lampinen and J. L. McClelland, "One-shot and few-shot learning of word embeddings," *arXiv preprint arXiv:1710.10280*, 2017.
- [338] C. Li, G. Wang, and G. De Melo, "Context-based few-shot word representation learning," in *Proc. IEEE Int. Conf. Semantic Computing (ICSC)*, 2018, pp. 239–242.
- [339] J. Sun, S. Wang, and C. Zong, "Memory, show the way: Memory based few shot word representation learning," in *Proc. Conf. Empirical Methods Nat. Language Process. (EMNLP)*, 2018, pp. 1435–1444.
- [340] Q. Liu, D. McCarthy, and A. Korhonen, "Second-order contexts from lexical substitutes for few-shot learning of word representations," in *Proc. Joint Conf. on Lexical and Computational Semantics*, 2019, pp. 61–67.
- [341] Z. Hu, T. Chen, K.-W. Chang, and Y. Sun, "Few-shot representation learning for out-of-vocabulary words," *arXiv preprint arXiv:1907.00505*, 2019.
- [342] M. Hofer, A. Kormilitzin, P. Goldberg, and A. Nevado-Holgado, "Few-shot learning for named entity recognition in medical text," *arXiv preprint arXiv:1811.05468*, 2018.
- [343] A. Fritzler, V. Logacheva, and M. Kretov, "Few-shot classification in named entity recognition task," in *Proc. ACM/SIGAPP Symposium on Applied Computing*, 2019, pp. 993–1000.
- [344] O. U. Florez and E. Mueller, "Learning to control latent representations for few-shot learning of named entities," *arXiv preprint arXiv:1911.08542*, 2019.
- [345] Z. Chen, H. Eavani, W. Chen, Y. Liu, and W. Y. Wang, "Few-shot nlg with pre-trained language model," *arXiv preprint arXiv:1904.09521*, 2019.
- [346] F. Mi, M. Huang, J. Zhang, and B. Faltings, "Meta-learning for low-resource natural language generation in task-oriented dialogue systems," in *Proc. Int. Joint Conf. Artificial Intell. (IJCAI)*. AAAI Press, 2019, pp. 3151–3157.
- [347] B. Peng, C. Zhu, C. Li, X. Li, J. Li, M. Zeng, and J. Gao, "Few-shot natural language generation for task-oriented dialog," *arXiv preprint arXiv:2002.12328*, 2020.
- [348] M. Kale and A. Rastogi, "Few-shot natural language generation by rewriting templates," *arXiv preprint arXiv:2004.15006*, 2020.
- [349] Y. Hou, Z. Zhou, Y. Liu, N. Wang, W. Che, H. Liu, and T. Liu, "Few-shot sequence labeling with label dependency transfer," *arXiv preprint arXiv:1906.08711*, 2019.
- [350] B. Lake, C.-y. Lee, J. Glass, and J. Tenenbaum, "One-shot learning of generative speech concepts," in *Proc. Annu. Meet. Cognit. Sci. Soc. (CogSci)*, vol. 36, no. 36, 2014.
- [351] J. Pons, J. Serrà, and X. Serra, "Training neural audio classifiers with few data," in *Proc. IEEE Int. Conf. Acoustics Speech Sig. Process. (ICASSP)*, 2019, pp. 16–20.
- [352] S. Zhang, Y. Qin, K. Sun, and Y. Lin, "Few-shot audio classification with attentional graph neural networks," *Proc. Interspeech*, pp. 3649–3653, 2019.
- [353] K.-H. Cheng, S.-Y. Chou, and Y.-H. Yang, "Multi-label few-shot learning for sound event recognition," in *Proc. IEEE Int. Workshop on Multimedia Sig. Process. (MMSP)*, 2019, pp. 1–5.
- [354] S.-Y. Chou, K.-H. Cheng, J.-S. R. Jang, and Y.-H. Yang, "Learning to match transient sound events using attentional similarity for few-shot sound recognition," in *Proc. IEEE Int. Conf. Acoustics Speech Sig. Process. (ICASSP)*, 2019, pp. 26–30.
- [355] J. Naranjo-Alcazar, S. Perez-Castanos, P. Zuccarello, and M. Cobos, "An open-set recognition and few-shot learning dataset

- for audio event classification in domestic environments," *arXiv preprint arXiv:2002.11561*, 2020.
- [356] Y. Chen, Y. Assael, B. Shillingford, D. Budden, S. Reed, H. Zen, Q. Wang, L. C. Cobo, A. Trask, B. Laurie *et al.*, "Sample efficient adaptive text-to-speech," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.
- [357] H. B. Moss, V. Aggarwal, N. Prateek, J. González, and R. Barra-Chicote, "Boffin tts: Few-shot speaker adaptation by bayesian optimization," in *Proc. IEEE Int. Conf. Acoustics Speech Sig. Process. (ICASSP)*, 2020, pp. 7639–7643.
- [358] Z. Zhang, Q. Tian, H. Lu, L.-H. Chen, and S. Liu, "Adadurian: Few-shot adaptation for neural text-to-speech with durian," *arXiv preprint arXiv:2005.05642*, 2020.
- [359] S. Choi, S. Han, D. Kim, and S. Ha, "Attentron: Few-shot text-to-speech utilizing attention-based variable-length embedding," *arXiv preprint arXiv:2005.08484*, 2020.
- [360] Y. Wang, J. Salamon, N. J. Bryan, and J. P. Bello, "Few-shot sound event detection," in *Proc. IEEE Int. Conf. Acoustics Speech Sig. Process. (ICASSP)*, 2020, pp. 81–85.
- [361] K. Shimada, Y. Koyama, and A. Inoue, "Metric learning with background noise class for few-shot detection of rare sound events," in *Proc. IEEE Int. Conf. Acoustics Speech Sig. Process. (ICASSP)*, 2020, pp. 616–620.
- [362] B. Shi, M. Sun, K. C. Puvvada, C.-C. Kao, S. Matsoukas, and C. Wang, "Few-shot acoustic event detection via meta learning," in *Proc. IEEE Int. Conf. Acoustics Speech Sig. Process. (ICASSP)*, 2020, pp. 76–80.
- [363] S. Arik, J. Chen, K. Peng, W. Ping, and Y. Zhou, "Neural voice cloning with a few samples," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 10019–10029.
- [364] B. Higy and P. Bell, "Few-shot learning with attention-based sequence-to-sequence models," *arXiv preprint arXiv:1811.03519*, 2018.
- [365] H. Seth, P. Kumar, and M. M. Srivastava, "Prototypical metric transfer learning for continuous speech keyword spotting with limited training data," in *Proc. Int. Workshop Soft Computing Models Industrial Environmental App.*, 2019, pp. 273–280.
- [366] D. Droghini, F. Vesperini, E. Principi, S. Squartini, and F. Piazza, "Few-shot siamese neural networks employing audio features for human-fall detection," in *Proc. Int. Conf. Pattern Recognit. Artificial Intell.*, 2018, pp. 63–69.
- [367] P. Anand, A. K. Singh, S. Srivastava, and B. Lall, "Few shot speaker recognition using deep neural networks," *arXiv preprint arXiv:1904.08775*, 2019.
- [368] Y. Duan, M. Andrychowicz, B. Stadie, O. J. Ho, J. Schneider, I. Sutskever, P. Abbeel, and W. Zaremba, "One-shot imitation learning," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 1087–1098.
- [369] C. Finn, T. Yu, T. Zhang, P. Abbeel, and S. Levine, "One-shot visual imitation learning via meta-learning," in *Proc. Conf. on Robot Learn.*, 2017, pp. 357–368.
- [370] T. Yu, C. Finn, A. Xie, S. Dasari, T. Zhang, P. Abbeel, and S. Levine, "One-shot imitation from observing humans via domain-adaptive meta-learning," *arXiv preprint arXiv:1802.01557*, 2018.
- [371] T. Yu, P. Abbeel, S. Levine, and C. Finn, "One-shot hierarchical imitation learning of compound visuomotor tasks," *arXiv preprint arXiv:1810.11043*, 2018.
- [372] Y. Aytar, T. Pfaff, D. Budden, T. Paine, Z. Wang, and N. de Freitas, "Playing hard exploration games by watching youtube," in *Proc. Advances Neural Inf. Process. Syst. (NeurIPS)*, 2018, pp. 2930–2941.
- [373] S. James, M. Bloesch, and A. J. Davison, "Task-embedded control networks for few-shot imitation learning," *arXiv preprint arXiv:1810.03237*, 2018.
- [374] D.-A. Huang, D. Xu, Y. Zhu, A. Garg, S. Savarese, L. Fei-Fei, and J. C. Niebles, "Continuous relaxation of symbolic planner for one-shot imitation learning," *arXiv preprint arXiv:1908.06769*, 2019.
- [375] Q. Shao, J. Qi, J. Ma, Y. Fang, W. Wang, and J. Hu, "Object detection-based one-shot imitation learning with an rgb-d camera," *Applied Sci.*, vol. 10, no. 3, p. 803, 2020.
- [376] A. Bonardi, S. James, and A. J. Davison, "Learning one-shot imitation from humans without humans," *IEEE Robotics Auto. Letters*, vol. 5, no. 2, pp. 3533–3539, 2020.
- [377] K. Frans, J. Ho, X. Chen, P. Abbeel, and J. Schulman, "Meta learning shared hierarchies," *arXiv preprint arXiv:1710.09767*, 2017.
- [378] A. Xie, A. Singh, S. Levine, and C. Finn, "Few-shot goal inference for visuomotor learning and planning," in *Proc. Conf. on Robot Learn.*, 2018, pp. 40–52.
- [379] D. Xu, S. Nair, Y. Zhu, J. Gao, A. Garg, L. Fei-Fei, and S. Savarese, "Neural task programming: Learning to generalize across hierarchical tasks," in *Proc. Int. Conf. Robotics Auto. (ICRA)*, 2018, pp. 1–8.
- [380] Y. Koizumi, S. Murata, N. Harada, S. Saito, and H. Uematsu, "Sniper: Few-shot learning for anomaly detection to minimize false-negative rate with ensured true-positive rate," in *Proc. IEEE Int. Conf. Acoustics Speech Sig. Process. (ICASSP)*, 2019, pp. 915–919.
- [381] A. Heidari, J. McGrath, I. F. Ilyas, and T. Rekatsinas, "Holodetect: Few-shot learning for error detection," in *Proc. Int. Conf. Management of Data*, 2019, pp. 829–846.
- [382] Y. Koizumi, M. Yasuda, S. Murata, S. Saito, H. Uematsu, and N. Harada, "Spidernet: Attention network for one-shot anomaly detection in sounds," in *Proc. IEEE Int. Conf. Acoustics Speech Sig. Process. (ICASSP)*, 2020, pp. 281–285.
- [383] X. Dong, L. Zhu, D. Zhang, Y. Yang, and F. Wu, "Fast parameter adaptation for few-shot image captioning and visual question answering," in *Proc. ACM Int. Conf. Multimedia (AMM)*, 2018, pp. 54–62.
- [384] D. Teney and A. van den Hengel, "Visual question answering as a meta learning task," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 219–235.
- [385] V. Prabhu, A. Kannan, M. Ravuri, M. Chaplain, D. Sontag, and X. Amatriain, "Few-shot learning for dermatological disease diagnosis," in *Mach. Learn. Healthcare Conf.*, 2019, pp. 532–552.
- [386] A. Paul, Y.-X. Tang, and R. M. Summers, "Fast few-shot transfer learning for disease identification from chest x-ray images using autoencoder ensemble," in *Medical Imaging 2020: Computer-Aided Diagnosis*, 2020.
- [387] W. Zhu, H. Liao, W. Li, W. Li, and J. Luo, "Alleviating the incompatibility between cross entropy loss and episode training for few-shot skin disease classification," *arXiv preprint arXiv:2004.09694*, 2020.
- [388] D. Rajan, J. J. Thiagarajan, A. Karayannidis, and S. Kashyap, "Self-training with improved regularization for few-shot chest x-ray classification," *arXiv preprint arXiv:2005.02231*, 2020.
- [389] S. Ali, B. Bhattacharya, T.-K. Kim, and J. Rittscher, "Additive angular margin for few shot learning to classify clinical endoscopy images," *arXiv preprint arXiv:2003.10033*, 2020.
- [390] C. Li, D. Zhang, Z. Tian, S. Du, and Y. Qu, "Few-shot learning with deformable convolution for multiscale lesion detection in mammography," *Medical Physics*, 2020.
- [391] X. Du, D. Zhong, and P. Li, "Low-shot palmprint recognition based on meta-siamese network," in *Proc. IEEE Int. Conf. Multimedia Expo (ICME)*, 2019, pp. 79–84.
- [392] H. Altae-Tran, B. Ramsundar, A. S. Pappu, and V. Pande, "Low data drug discovery with one-shot learning," *ACS Central Sci.*, vol. 3, no. 4, pp. 283–293, 2017.
- [393] J. Liu, S. J. Gibson, J. Mills, and M. Osadchy, "Dynamic spectrum matching with one-shot learning," *Chemometrics Intell. Laboratory Syst.*, vol. 184, pp. 175–181, 2019.
- [394] Y. Li and J. Yang, "Few-shot cotton pest recognition and terminal realization," *Computers and Electronics in Agriculture*, vol. 169, pp. 1–9, 2020.
- [395] M. M. U. Chowdhury, F. Hammond, G. Konowicz, C. Xin, H. Wu, and J. Li, "A few-shot deep learning approach for improved intrusion detection," in *Proc. IEEE Annu. Ubiquitous Computing, Electronics and Mobile Communication Conf. (UEMCON)*, 2017, pp. 456–462.
- [396] T. Gong, Y. Kim, J. Shin, and S.-J. Lee, "Metasense: few-shot adaptation to untrained conditions in deep mobile sensing," in *Proc. Conf. Embedded Networked Sensor Syst.*, 2019, pp. 110–123.
- [397] K.-S. Fu and J. Mui, "A survey on image segmentation," *Pattern Recognit.*, vol. 13, no. 1, pp. 3–16, 1981.
- [398] L. Liu, W. Ouyang, X. Wang, P. Fieguth, J. Chen, X. Liu, and M. Pietikäinen, "Deep learning for generic object detection: A survey," *Int. J. Comput. Vis.*, vol. 128, no. 2, pp. 261–318, 2020.
- [399] D. Brezale and D. J. Cook, "Automatic video classification: A survey of the literature," *IEEE Trans. on Sys., Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 38, no. 3, pp. 416–430, 2008.
- [400] Z. Shou, D. Wang, and S.-F. Chang, "Temporal action localization in untrimmed videos via multi-stage cnns," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 1049–1058.
- [401] R. Yao, G. Lin, S. Xia, J. Zhao, and Y. Zhou, "Video object segmentation and tracking: A survey," *arXiv preprint arXiv:1904.09172*, 2019.

- [402] V. Korde and C. N. Mahender, "Text classification and classifiers: A survey," *Int. J. Artificial Intell. App.*, vol. 3, no. 2, p. 85, 2012.
- [403] J. Camacho-Collados and M. T. Pilehvar, "From word to sense embeddings: A survey on vector representations of meaning," *J. Artificial Intell. Research*, vol. 63, pp. 743–788, 2018.
- [404] M. Nickel, K. Murphy, V. Tresp, and E. Gabrilovich, "A review of relational machine learning for knowledge graphs," *Proc. of the IEEE*, vol. 104, no. 1, pp. 11–33, 2015.
- [405] J. Kober, J. A. Bagnell, and J. Peters, "Reinforcement learning in robotics: A survey," *Int. J. Robotics Research*, vol. 32, no. 11, pp. 1238–1274, 2013.
- [406] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne, "Imitation learning: A survey of learning methods," *ACM Computing Surveys*, vol. 50, no. 2, pp. 1–35, 2017.
- [407] F. Bonin-Font, A. Ortiz, and G. Oliver, "Visual navigation for mobile robots: A survey," *J. Intell. Robotic Syst.*, vol. 53, no. 3, p. 263, 2008.
- [408] M. P. Deisenroth, G. Neumann, and J. Peters, *A survey on policy search for robotics*. now publishers, 2013.
- [409] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," *ACM Computing Surveys*, vol. 41, no. 3, pp. 1–58, 2009.
- [410] J. Lu, S. Jin, J. Liang, and C. Zhang, "Robust few-shot learning for user-provided data," *IEEE Trans. Neural Net. Learn. Syst.*, 2020.
- [411] Y.-J. Park and A. Tuzhilin, "The long tail of recommender systems and how to leverage it," in *Proc. ACM Conf. Recommender Syst.*, 2008, pp. 11–18.
- [412] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba, "Network dissection: Quantifying interpretability of deep visual representations," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 6541–6549.
- [413] Q.-s. Zhang and S.-C. Zhu, "Visual interpretability for deep learning: a survey," *Frontiers Infor. Tech. Electronic Eng.*, vol. 19, no. 1, pp. 27–39, 2018.



Jiang Lu received the B.S. and Ph.D. degree from the Department of Automation, Tsinghua University, Beijing, China, in 2013 and 2020, respectively. He is currently an Associate Research Assistant with China Marine Development and Research Center (CMDRC), Beijing, China.

He has served as a Reviewer for the IEEE Transactions on Pattern Analysis and Machine Intelligence, and Pattern Recognition. His research interests include machine learning, deep learning and computer vision.

- [414] J. Lu, L. Li, and C. Zhang, "Self-reinforcing unsupervised matching," *arXiv preprint arXiv:1909.04138*, 2019.



Jieping Ye (Fellow, IEEE) received the Ph.D. degree in computer science from the University of Minnesota, Twin Cities, MN, USA, in 2005.

He is the Head of Didi AI Labs, VP of Didi Chuxing, and a Didi Fellow. He is also a Professor with the University of Michigan, Ann Arbor, MI, USA. His research interests include big data, machine learning, and data mining with applications in transportation and bio-medicine. He was a recipient of the NSF CAREER Award in 2010. His papers have been selected for the Outstanding

Student Paper at ICML in 2004, the KDD Best Research Paper Runner Up in 2013, and the KDD Best Student Paper Award in 2014. He has served as a Senior Program Committee/Area Chair/Program Committee Vice Chair of many conferences, including NIPS, ICML, KDD, IJCAI, ICDM, and SDM. He serves as an Associate Editor for Data Mining and Knowledge Discovery, the IEEE Transactions on Knowledge and Data Engineering, and the IEEE Transactions on Pattern Analysis and Machine Intelligence.



Jianwei Zhang (M'95) received the B.Eng. (Hons.) and M.Eng. degrees from the Department of Computer Science, Tsinghua University, Beijing, China, in 1986 and 1989, respectively, the Ph.D. degree from the Institute of Real-Time Computer Systems and Robotics, Department of Computer Science, University of Karlsruhe, Karlsruhe, Germany, in 1994, and the Habilitation from the Faculty of Technology, University of Bielefeld, Bielefeld, Germany, in 2000.

He is currently a Professor and the Head of the TAMS Group, Department of Informatics, University of Hamburg, Hamburg, Germany. He has published about 300 journal and conference papers (winning four best paper awards), technical reports, four book chapters, and five research monographs. He has been coordinating numerous collaborative research projects of EU and German Research Council, including the Transregio-SFB TRR 169 "Crossmodal Learning." His current research interests include cognitive robotics, sensor fusion, dexterous manipulation, and multimodal robot learning. Dr. Zhang is a Life-Long Academician of the Academy of Sciences, Hamburg. He is the General Chair of the IEEE MFI 2012, the IEEE/RSJ IROS 2015, and the IEEE Robotics and Automation Society AdCom from 2013 to 2015.



Pinghua Gong received the B.S. degree from the Department of Automation, Xi'an Jiaotong University, Xi'an, China, in 2008, and the Ph.D. degree from the Department of Automation, Tsinghua University, Beijing, China, in 2013.

He is currently a Distinguished Scientist & Senior Director at Didi Chuxing, Beijing. He has published more than 20 papers in top-tier journals and conferences, including JMLR, NIPS, ICML, KDD, IJCAI, and AAAI. His research interests include machine learning and data mining.



Changshui Zhang (Fellow, IEEE) received the B.S. degree in mathematics from Peking University, Beijing, China, in 1986, and the M.S. and Ph.D. degrees in control science and engineering from Tsinghua University, Beijing, in 1989 and 1992, respectively.

He is currently a Professor with the Department of Automation, Tsinghua University. His research interests include artificial intelligence, image processing and machine learning.