

Review Article

The Evolution of Knowledge Distillation in Image Classification Tasks

Tiejun Yin* 

School of Artificial Intelligence, Nanjing Normal University of Special Education, Nanjing, China

Abstract

With the start of the digital age, image classification has been one of the first problems in computer vision, that is to say, to identify and label all pictures according to various categories. High-performance models are usually very expensive to run on computers, so they are not suitable for the resource-constrained environment of real-time deployment. One of the many methods to address the problem of model compression is knowledge distillation. The goal is to reduce the size and complexity of the "teacher model" by optimising the parameters of the "student model" so that it can achieve similar performance to the "teacher model", but with a lower cost and without modifying the "student model" architecture. In this paper, we aim to systematically review the technical methods of knowledge distillation in image classification and discuss how to maximise the learning efficiency of student models with different forms of knowledge transfer across various dimensions, such as model output and feature map matching, structural relationships between the teacher and student models, etc. The three parts of the study are as follows: Output-based distillation is how to separate knowledge distillation from the problem of determining the boundary of a relationship graph; Feature-based distillation is related to Wasserstein distance and feature attention transfer; and Relationship-based distillation is based on virtual distillation techniques and Lipschitz continuity constraints. Based on the above analysis, although there is a problem of one-dimensional transfer, complementary effects can be obtained through the organic combination of output, feature and relationship transfer. A multi-dimensional fusion method can improve the accuracy and generalisation of the small model so that it does not need to be implemented on the high-performance vision system of the edge computing platform.

Keywords

Knowledge Distillation, Computer Vision, Model Compression

1. Introduction

The technological advancements and development of deep neural networks in the field of artificial intelligence have ushered in a historic transformation in computer vision. Deep learning acts as the driver in image classification. Compared with other models, Convolutional Neural Network (CNN) uses convolutional kernels for local perception, which helps

to reduce the amount of parameters and make the model more efficient. They extract features through multiple layers of convolution and pooling.

Application: Image classification is applied in the medical field, which includes medical imaging diagnostics based on X-ray images, CT images and MRI images. In the security

*Correspondence: Tiejun Yin (073arfstu@stu.njts.edu.cn)

Received: 28 May 2026; Accepted: 10 June 2026; Published: 26 June 2026



field, such as facial recognition, dynamic monitoring, behavior scanning and analysis, etc. However, image classification tasks are constrained by data resources and model quality, and dataset imbalance also impacts the model's generalization capability. Deep neural networks are often regarded as black boxes, limiting model reliability and rendering deep learning computations too complex to meet real-time requirements. The introduction of knowledge distillation effectively addresses the challenge of deploying overly large datasets required for training large models. By transferring knowledge from complex, massive yet high-performance teacher models to student models, it reduces the performance gap between them.

This paper employs the methodology proposed by classical knowledge distillation to progressively deepen the analysis to output-based structural, feature-level, and relational theoretical dimensions, providing a comprehensive overview of the primary strategies and key frameworks within knowledge distillation for the task of image classification.

2. Output-Oriented Knowledge Distillation

Knowledge Distillation (KD) is an outcome-oriented model compression technique originally proposed by Hinton et al. for image classification. It improves the performance of a small model (student model) by training it with the help of a large model (teacher model) [1]. This knowledge initially refers to the category probabilities obtained from the softmax layer of the teacher model. The student model takes the category probability vector as its learning target and mimics the teacher model's classification behavior. Distillation refers to raising the parameter temperature in the softmax layer to soften the output probability vector, making it easier for the student model to learn and thereby improving the generalization performance of the model. By improving the output-based knowledge distillation method and dividing the knowledge distillation into two dimensions, network performance and network compression, the method has improved logical coherence and interpretability [2]. Reo Fukunaga et al. proposed a closed-loop ACD-U (Asymmetric Co-Teaching with Machine Unlearning) framework to tackle the confirmation bias problem and memory accumulation issue in noisy label learning. This framework comprises an early stage stabilising model and a convolutional neural network, to create a pair of complementary learners: (1) one is used to suppress the noisy model through a Gaussian mixture model and confidence thresholds, and (2) one is used to dynamically identify and forget the inconsistent memories with the loss trajectory and image-language training. This innovative mechanism is able to overcome the limitation of 'One Way Defense' and achieves outstanding performance on standard benchmarks like CIFAR-100. [3].

2.1. Structural Innovation

As classical knowledge distillation is knowledge transferring based on the category probabilities output by the softmax layer, Borui Zhao et al. revisited the classical KD. Upon classical KD, they classified this classical KD into two approaches: Target-Class Knowledge Distillation (TCKD) aims at answering the target-class macro-level judgment question, i.e., whether the category is the target class [1]; and Non-Target-Class Knowledge Distillation (NCKD) focuses on how to suppress the non-target classes [4]. Specifically, the TCKD transmits prediction information for target categories, providing sample difficulty signals in the form of a binary classification task, while the NCKD reflects the similarity structure among negative categories. Due to coupling issues in the classic KD loss, which suppress NCKD effectiveness and prevent independent adjustment of TCKD and NCKD weights, Decoupled Knowledge Distillation (DKD) is proposed to reduce coupling between the two and enhance efficiency. Based on classical knowledge distillation, DKD redefines the loss function as shown in Equation (1). Since the weights of NCKD are influenced by the target category probabilities of the teacher model, hyperparameters α and β are introduced to control the TCKD and NCKD components respectively, enabling truly independent adjustment of each part.

$$DKD = \alpha TCKD + \beta NCKD \quad (1)$$

2.2. The Essence of Distillation

Although researchers have explored various knowledge distillation techniques and their applications, such as Son, W et al. proposing a densely guided knowledge distillation technique that introduces multiple intermediate-scale teacher assistants for multi-stage knowledge distillation, enabling knowledge transfer between large-scale teacher models and small-capacity student models [5]. However, the understanding of the essence of knowledge distillation remains incomplete. To address this, Utkarsh Ojha et al. proposed a geometric interpretation framework based on decision boundaries. For the first time, they analyzed that latent knowledge is not only the inter-class relationships provided by soft labels but also a manifestation of the teacher's decision-making system. This confirmed that knowledge distillation essentially involves the student model learning the decision-making mechanism of the teacher model, thereby holistically shaping the student's decision-making behavior. It has been demonstrated that knowledge distillation also conveys implicit properties such as robustness, data invariance, and color constancy [6].

3. Feature-Based Knowledge Distillation

The classic KD method focuses solely on the final result, neglecting the characteristics of the teacher network's intermediate layers. This makes it increasingly difficult to optimize

deeper networks as shallower student models mimic the teacher model's structure. To this end, Romero et al. utilized feature knowledge from intermediate layers to guide student network training. They selected the intermediate layer of the teacher network as the prompt layer and the intermediate layer of the student network as the guidance layer. By adding a convolutional regressor on the guidance layer, they minimized the loss between the two networks. They adopted a strategy of training the student model with fixed teacher model parameters from the input to the guidance layer, followed by global network knowledge distillation training. They pioneered the introduction of prompt learning mechanisms to train deep models, breaking through the limitations of traditional output knowledge distillation [7]. Ziyao Guo et al. proposed a Category-Attention-Transferring Knowledge Distillation (CAT-KD) with strong interpretability. They demonstrated that transferring only the category activation map enhances the student model's ability to distinguish strong discriminative regions and guides it to focus on more important areas. The overall loss is the sum of cross-entropy loss and CAT loss with an introduced β balancing factor [8]. Liuchi Xu et al. introduced a new approach called Heterogeneous Complementary Distillation (HCD) to tackle the feature representation mismatch problem due to bias differences between knowledge distillation with different architectures. The shared logits are decomposed and fused with teacher logits according to the low-level information and high-level semantic information using a Complementary Feature Mapper (CFM), so as to improve classification consistency; Orthogonal Loss is also introduced to promote the diversity of decomposed sub-logits and prevent redundant knowledge transfer. The experiments eventually became successful, effectively combining the benefit of global modelling of the teacher and local feature modelling of the student, which greatly improves the stability of the model and generalisation performance. [9].

Although The model proposed by Romero et al. introduces a novel perspective by aligning intermediate-layer features between teacher and student networks to enhance student model performance, shallow feature alignment may lead to semantic inconsistencies when model structures differ significantly. As feature extraction grows more complex, the required transformation functions and alignment strategies also become more intricate, and processing higher-dimensional feature information demands greater computational resources.

Since the Leibler Divergence (KL-Div) in the same category has achieved certain applications in many fields and gained good performance, there are also certain limitations in KL-Div. KL-Div is only used for the same category comparison, can't be extended to inter-category comparison directly, and has certain problems when there is no overlapping region in the intermediate feature comparison. So Jiaming Lv et al. designed a knowledge distillation method based on Wasserstein Distance (WD) to compete with KL-Div, which can be divided into discrete Logits Distillation (WKD-L) and continuous Feature Distillation (WKD-F) [10].

The paper defines the WD as the minimum cost in. WKD-L employs a discrete WD to measure prediction discrepancies between models, quantifying category relationships through Centered Kernel Alignment (CKA) to enable comparisons across different categories [11]. WKD-F models intermediate layer features using Gaussian distributions, matching feature maps via continuous WD while considering the Riemannian manifold geometry of positive definite symmetric matrices to effectively transfer deep feature knowledge. While WD-based knowledge distillation holds immense potential to overcome traditional method limitations, it still faces challenges such as high computational overhead and feature modeling constraints.

4. Relational Knowledge Distillation

Most feature-based and output-based knowledge distillation methods focus solely on knowledge within independent samples, whereas relational knowledge distillation places greater emphasis on structural knowledge within models and deep exploration of categorical relationships. Chuanguang Yang et al. discussed the general form of distillation losses as shown in Equation (2):

$$\mathcal{L}_{relation_{kd}}(F^S, F^T) = \sum_{ij} \mathcal{L}_{dis}(\psi^S(v_i^S, v_j^S), \psi^T(v_i^T, v_j^T)) \quad (2)$$

F^S, F^T representing the feature sets of the student model and teacher model respectively. v_i, v_j denote the feature embeddings of the i -th and j -th samples, respectively. ψ^S and ψ^T represent similarity measures for sample feature embeddings, \mathcal{L}_{dis} while serves as the distance function for instance graph similarity [12]. Hu Chengming et al. have abandoned the original teacher-student system, and rebuilt the system into a general purpose knowledge transfer operating system, which is no longer limited to the single compression paradigm, and introduced a complete theoretical system, including compression and expansion objectives. The paper presents the first end-to-end framework that, besides making it clearer what has been common to distillation, also provides some directions such as architectural coordination, knowledge quality, etc., so that knowledge distillation can go from a technique/approach in the realm of engineering to become a systematic field in the realm of knowledge engineering. [13].

Traditional relational knowledge distillation is weak in inducing relational matching, which leads to overfitting and interference from false information. Thus, its performance is far below instance-matching method. To conquer these problems, Weijia Zhan et al. proposed a new architecture, Virtual Relational Matching Knowledge Distillation (VRM) to help student model learn more informative affinity graphs which are rich in sample information, inter-class relations and inter-view structural relationships [14].

Firstly, they employ dense relationship graphs to learn inter-sample relationships from predicted logits. Then, they design a category-batch-level relationship graph structure to preserve response variations to learn structural knowledge. After

that, they develop virtual graphs to learn virtual-real relationships. What is more, to alleviate false gradients, the affinity graph will be pruned twice to remove redundant edges at both source and target sides. At last, they combine huber loss with cross entropy loss.

A Novel Knowledge Distillation Approach learns virtual knowledge to improve model training revisits classical relational distillation and makes new improvements. For the first time, virtual relationship is introduced into graph structure of knowledge distillation. The exploration of relational distillation is reactivated and a larger search space is developed.

But the traditional feature based knowledge distillation only align the shallow-level knowledge, treating neural networks as black box and ignoring higher-level knowledge, i.e, functional features, which leads to the student directly imitating the teacher in a simple way. Thus, to bridge the gap, Shang Y. et al. proposed Lipschitz-Guided Knowledge Distillation (LONDON). They utilize Lipschitz continuity as the knowledge to be transferred and obtain the knowledge transfer by minimizing the distance between Lipschitz constants of teacher-student networks [15].

Lipschitz continuity is typically defined as follows: for a

real-valued function $f: X \rightarrow Y$, where X and Y are metric spaces, if there exists a constant $L \geq 0$ such that for all $x_1, x_2 \in X$, the difference between the two functions, i.e., $|f(x_1) - f(x_2)|$, and the rate of change between any two points is $\leq L$, then L is called the Lipschitz constant. Although computing the Lipschitz constant is prohibitively difficult, the paper proposes approximating it using the transfer matrix of independent modules. If this matrix is normalized to become orthogonal, the spectral norm of the weight matrix can be obtained by calculating the maximum eigenvalue of the transfer matrix (avoiding direct computation of large-scale matrices). This approach approximates the Lipschitz constant for each module and employs a power iteration method for global network approximation.

The proposed LONDON knowledge distillation breakthrough the aforementioned limitations on solely considering shallow knowledge between features and outputs. Efficient approximation of constants through the power iteration method of transfer matrices, that serves as a solid theoretical support and available extension for relational distillation.

5. Comparative Analysis and Discussion

Table 1. Comparison of Pros and Cons of Knowledge Distillation.

	Output-Oriented Knowledge Distillation	Feature-Based Knowledge Distillation	Relational Knowledge Distillation
Advantages	Easy to implement and highly scalable, suitable for multiple tasks such as classification and detection.	Capable of capturing abstract details and conveying richer semantic information	Uncover deeper relationships, focus on model structural characteristics, and demonstrate strong generalization capabilities.
Disadvantages	Focusing solely on outcomes while neglecting feature information imposes significant limitations.	Cannot be applied to models with excessive structural differences; simple alignment operations introduce noise.	High implementation difficulty, significant computational resource consumption, unstable distillation results

Based on learning distillation and comparison in strengths and weaknesses (Table 1), knowledge network strategies to improve network effectiveness by exploring and innovating continuously have achieved remarkable performance as shown in Table 1.

Network knowledge ways from output distillation (a performance-centered approach which employs soft labels to express implicit knowledge), to feature distillation (which expresses abstract feature information for more concise performance compression), to relationship distillation (which establishes complicated relationships with models).

However, there are still many issues existing in knowledge distillation research, such as how to design effective metrics to evaluate knowledge distillation or not, how to implement

tasks efficiently with huge and complicated data in ultra-high-precision scene, how to evaluate the transparency and regularity of knowledge distillation process, and how to solve the poor interpretability issue of some implicit information which may suppress distillation effectiveness. So, in the future, people should consider using more kinds of knowledge sources to improve the controllability of the model and satisfy the requirement of legitimacy and privacy. Unlike traditional transmission ways, students can learn how to communicate logically from teacher’s network. Finally, people make a breakthrough in knowledge distillation in the view of green AI and sustainability.

6. Conclusions

This paper discusses the knowledge distillation methods for image classification in three different levels: output-based, feature-based and relationship-based. The simplest approach, which trains small models to mimic the outputs of large models, is straightforward and inexpensive but has inherent limitations.

Feature-based knowledge distillation uses feature knowledge for distillation, while only shallow-level feature information is aligned.

Unlike the previous studies based on single-sample, Relational knowledge distillation emboldens to propose distillation based on sample relationships. Theoretical refinement leads to a new attempt for knowledge transfer with higher-level knowledge---function characteristics. Finally, the novel structured knowledge distillation effectively overcomes the weakness of classical knowledge distillation.

In contrast to classical knowledge distillation, in the future, the black-box models will be gradually tuned to be white-box models. Federated knowledge distillation, cross-modal joint distillation, multi-teacher knowledge fusion and PEFT technologies (including LoRA, Adapter etc.) will be widely used.

Abbreviations

CNN	Convolutional Neural Network
KD	Knowledge Distillation
TCKD	Target-Class Knowledge Distillation
NCKD	Non-Target-Class Knowledge Distillation
DKD	Decoupled Knowledge Distillation
KL-Div	Leibler Divergence
WD	Wasserstein Distance
WKD-L	Logits Distillation
WKD-F	Feature Distillation
CKA	Centered Kernel Alignment
VRM	Virtual Relational Matching Knowledge Distillation
LONDON	Lipschitz-Guided Knowledge Distillation
HCD	Heterogeneous Complementary Distillation
CFM	Complementary Feature Mapper
ACD-U	Asymmetric Co-Teaching with Machine Unlearning

Author Contributions

Tiejun Yin: Conceptualization, Investigation, Data curation, Writing – original draft, Writing – review & editing

Conflicts of Interest

The author declares no conflicts of interest.

References

- [1] Hinton G, Vinyals O, Dean J. Distilling the Knowledge in a Neural Network. 2015. <https://doi.org/10.48550/arXiv.1503.02531>
- [2] Qi Honggang, Si Zhaofeng. A Review of the Research and Application of Knowledge Distillation Methods. *Journal of Image and Graphics*, 2023, 28(09): 2817-2832. <https://dx.doi.org/10.11834/jig.220273>
- [3] Fukunaga R, Yoshida S, Muneyasu M. ACD-U: Asymmetric co-teaching with machine unlearning for robust learning with noisy labels. 2026. <https://doi.org/10.48550/arXiv.2603.07166>
- [4] Zhao B, Cui Q, Song R, Qiu Y, Liang J. Decoupled Knowledge Distillation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022: 11953-11962. <https://doi.org/10.48550/arXiv.2203.08679>
- [5] Son W, Na J, Choi J, Hwang W. Densely Guided Knowledge Distillation Using Multiple Teacher Assistants. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021: 9395-9404. <https://doi.org/10.48550/arXiv.2009.08825>
- [6] Ojha U, Li Y, Sundara Rajan A, Liang Y, Lee Y J. What Knowledge Gets Distilled in Knowledge Distillation? *Advances in Neural Information Processing Systems*, 2023, 36: 11037-11048. <https://doi.org/10.48550/arXiv.2205.16004>
- [7] Romero A, Ballas N, Kahou S E, Chassang A, Gatta C, Bengio Y. FitNets: Hints for Thin Deep Nets. *arXiv Preprint arXiv: 1412.6550*, 2014. <https://doi.org/10.48550/arXiv.1412.6550>
- [8] Guo Z, Yan H, Li H, Lin X. Class Attention Transfer Based Knowledge Distillation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023: 11868-11877. <https://doi.org/10.48550/arXiv.2304.12777>
- [9] XU L, LIU K, LIU J, et al. Heterogeneous Complementary Distillation [C]// *Proceedings of the Thirty-Ninth AAAI Conference on Artificial Intelligence (AAAI-26)*. Vancouver, BC, Canada: AAAI Press, 2025: 11316–11322. <https://doi.org/10.1609/aaai.v40i13.38112>
- [10] Lv J, Yang H, Li P. Wasserstein Distance Rivals Kullback-Leibler Divergence for Knowledge Distillation. *Advances in Neural Information Processing Systems*, 2024, 37: 65445-65475. <https://doi.org/10.48550/arXiv.2412.08139>
- [11] Cortes C, Mohri M, Rostamizadeh A. Algorithms for Learning Kernels Based on Centered Alignment. *Journal of Machine Learning Research*, 2012, 13(1): 795-828. <https://doi.org/10.48550/arXiv.1203.0550>
- [12] Yang C, Yu X, An Z, Xu Y. Categories of Response-Based, Feature-Based, and Relation-Based Knowledge Distillation. In: *Advancements in Knowledge Distillation: Towards New Horizons of Intelligent Systems*. Cham: Springer International Publishing, 2023: 1-32. <https://doi.org/10.48550/arXiv.2306.10687>
- [13] HU C, LI X, LIU D, et al. Teacher-Student Architecture for Knowledge Distillation: A Survey [EB/OL]. 2023. <https://doi.org/10.48550/arXiv.2308.04268>

- [14] Zhang W, Xie F, Cai W, Ma C. VRM: Knowledge Distillation via Virtual Relation Matching. arXiv Preprint, 2025. <https://doi.org/10.48550/arXiv.2502.20760>
- [15] Shang Y, Duan B, Zong Z, Nie L, Yan Y. Lipschitz Continuity Guided Knowledge Distillation. Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021: 10675-10684. <https://doi.org/10.48550/arXiv.2108.12905>