Supplementary: Deep ensemble learning by diverse knowledge distillation for fine-grained object classification

Naoki Okamoto, Tsubasa Hirakawa, Takayoshi Yamashita, and Hironobu Fujiyoshi

Chubu University, Kasugai, Aichi, Japan {naok,hirakawa}@mprg.cs.chubu.ac.jp,{takayoshi,fujiyoshi}@isc.chubu.ac.jp

1 Comparison of loss designs

In this section, we show the effect that loss designs have on ensembles in terms of bringing knowledge closer and separating knowledge. We trained four networks and compared the ensemble accuracy using them. We used ResNet-18 [2] as the network, the probability distribution and attention map as the knowledge, and Stanford Dogs [3] as the dataset. The attention map was created from the output of ResBlock4 using attention transfer [8]. The loss design for the probability distribution used KL-divergence (KL) and cosine similarity (cos). The loss design for the attention map used mean square error (MSE) and cosine similarity (cos).

Tables 1 and 2 show the results of the loss design for the probability distribution and the attention map. "+" means the loss design for bringing knowledge closer and "-" means the loss design for separating knowledge. With the exception of several loss designs, the ensemble accuracy did not change significantly with the loss design. Loss designs that use division have the possibility of dividing by zero. Therefore, to train the network as a minimization problem, we selected different loss designs for bringing knowledge closer and separating knowledge.

2 Types of knowledge and effects of loss design

In this section, we show the effect of knowledge type and loss design on network accuracy. In Sec. 3.2 of the main paper, the ensemble accuracy (Fig. 3 in the

Table 1: Comparison of ensemble accuracy for different loss designs with probability distribution [%].

Loss design	Knowledge	Ensemble accuracy
KL	+	69.95
$-\cos$	+	70.00
$1/\cos$	+	69.12
cos	—	69.88
-KL	-	57.87
1/KL	-	67.88

Table 2: Comparison of ensemble accuracy for different loss designs with attention map [%].

Loss design	Knowledge	Ensemble accuracy
MSE	+	68.91
$-\cos$	+	68.93
$1/\cos$	+	68.98
cos	-	70.35
-MSE	-	64.73
1/MSE	—	68.50

2 N. Okamoto et al.



Fig. 1: Network accuracy by diverse knowledge distillation [%].

main paper) is comparable depending on the loss design. The accuracy of the network in Sec. 3.2 is shown in Fig. 1. The accuracy of the network tends to vary depending on the loss design. Therefore, we think that different loss designs have different learning effects. Separating the probability distributions causes a decrease in the accuracy of the network, whereas separating the attention maps prevents a decrease in the accuracy of the network.

3 Process of graph optimization

Algorithm 1 shows process of graph optimization. Each graph is evaluated at the timing of 2^k epochs. If the ensemble accuracy of the graph is higher than the median ensemble accuracy of previously evaluated graphs at the same epoch timing, training is continued to the next 2^k epochs and evaluated again. If the ensemble accuracy of the graph is lower than the median, the training of the graph is terminated. Then, a new graph is created for the structure that has not yet been evaluated, and training of the new graph is started from 1 epoch. These processes are repeated until the number of graphs created reaches 6,000.

4 Visualization of optimized graphs

In the main paper, we showed a graph structure that was automatically designed using Stanford Dogs. In this section, we show graph structures that were automatically designed using datasets other than Stanford Dogs. Figures 2, 3, 4, and 5 show graphs automatically designed using CIFAR-10 [5], CIFAR-100 [5], Caltech-UCSD Birds-200-2011 (CUB-200-2011) [7], and Stanford Cars [4]. A red node represents an ensemble node, a gray node represents a network node, and "Label" represents supervised labels. At each edge, the selected loss design and gate are shown, excluding the cutoff gate. The accuracy in parentheses is the result for the dataset used for automatic graph design. Deep ensemble learning by diverse KD for fine-grained object classification

Algorithm 1 Optimizing the graph of *M* network nodes

Require: Number of searches N_{search} , Number of GPUs N_{gpu} **Require:** Training epochs T, Train set D_{train} , Test set D_{test} **Require:** Networks set $\mathcal{F} : \{f_1, ..., f_M\}$, Search space of graphs $\mathcal{S} : \{S_1, ..., S_i\}$ **Require:** Random sampling function ϕ , Pruning function ASHA, The function computing ensemble accuracy ψ 1: $A_{\text{best}} \leftarrow 0$ 2: $N_{\text{end}} \leftarrow 0$ 3: Start the asynchronous search in $N_{\rm gpu}$ environments 4: while $N_{\text{end}} \ge N_{\text{search}} \mathbf{do}$ Initialization of networks \mathcal{F} 5: $S \leftarrow \phi(\mathcal{S})$ 6: // Selection of graph structure by random sampling 7:for $t \leftarrow 1$ to T do 8: $S(\mathcal{F}, D_{\text{train}})$ // Training of networks by graph // Evaluation of ensemble accuracy 9: $A_{\text{ens}} \leftarrow \psi(\mathcal{F}, D_{\text{test}})$ 10: // Judgement of pruning by ASHA 11: if $ASHA(A_{ens})$ then break 12:end if 13:14:end for 15: $A_{\text{ens}} \leftarrow \psi(\mathcal{F}, D_{\text{test}})$ 16:if $A_{ens} > A_{best}$ then 17: $S_{\text{best}} \leftarrow S$ 18: $A_{\text{best}} \leftarrow A_{\text{ens}}$ end if 19: $N_{\text{end}} \leftarrow N_{\text{end}} + 1$ 20:21: end while 22: return S_{best}





Fig. 2: Two-node graph optimized on CIFAR-10.

Fig. 3: Two-node graph optimized on CIFAR-100.

4 N. Okamoto et al.



Fig. 4: Graph optimized on CUB-200-2011.



Fig. 5: Graph optimized on Stanford Cars.

Deep ensemble learning by diverse KD for fine-grained object classification

5 Generalizability of graphs

In the main paper, we showed the generalizability of graph structures automatically designed using Stanford Dogs for four datasets. In this section, we evaluate five graphs automatically designed using five different datasets. We used attention branch network (ABN) [1] based on ResNet as the network. We used Stanford Dogs, Stanford Cars, CUB-200-2011, CIFAR-10, and CIFAR-100 as the datasets. Stanford Dogs, Stanford Cars, and CUB-200-2011 belong to the fine-grained object classification task. CIFAR-10 and CIFAR-100 belong to the general object classification task.

Table 3 shows the ensemble accuracy of the two-node graph. "Independent" is the result of an individually trained network. "DML" is the result of a network trained with deep mutual learning [9]. The bold text in the Ensemble column shows that the ensemble accuracy was higher than Independent and DML. The graphs automatically designed using the datasets of the fine-grained object classification task showed an improved ensemble accuracy, especially in the fine-grained object classification task. The graphs automatically designed using the datasets of the general object classification task showed an improved ensemble accuracy, especially in the general object classification task. We believe that there was generalizability in the graph structure when the problem set was the same and that optimization resulted in a graph structure that corresponded to the problem set.

Tables 4, 5, and 6 show the results for various numbers of network nodes, from two to five, between the fine-grained object classification tasks. We believe that there was generalizability in the graph structure even when the number of network nodes was increased. Focusing on the ensemble accuracy of "Ours," the ensemble accuracy varied depending on the dataset used for the automatic design. We believe that this is because the number of combinations of graph structures was huge, and automatic design using random search finally resulted in different graph structures.

Figures 6a and 6b show the attention maps of the five-node graph optimized by Stanford Dogs for training CUB-200-2011 and Stanford Cars, respectively. Each point in the graph is two-dimensional because the attention maps were reduced to two dimensions by UMAP [6]. We see that the attention maps of each dataset are similar even when training on a dataset different from the optimization. These dimensionally reduced maps have a similar trend for each node regardless of the sample. This indicates that the graph structure is generalizable in terms of the attention map.

Method	Training	Optimizing	Ensemble	
Method	Graph	Graph		
Independent	Stanford Dogs	-	70.90	
DML	Stanford Dogs	-	71.45	
Ours	Stanford Dogs	Stanford Dogs	73.86	
Ours	Stanford Dogs	CUB-200-2011	72.43	
Ours	Stanford Dogs	Stanford Cars	72.79	
Ours	Stanford Dogs	CIFAR-100	71.53	
Ours	Stanford Dogs	CIFAR-10	70.08	
Independent	CUB-200-2011	-	65.26	
DML	CUB-200-2011	-	66.90	
Ours	CUB-200-2011	Stanford Dogs	72.06	
Ours	CUB-200-2011	CUB-200-2011	69.81	
Ours	CUB-200-2011	Stanford Cars	71.27	
Ours	CUB-200-2011	CIFAR-100	66.43	
Ours	CUB-200-2011	CIFAR-10	66.42	
Independent	Stanford Cars	-	88.49	
DML	Stanford Cars	-	88.89	
Ours	Stanford Cars	Stanford Dogs	89.76	
Ours	Stanford Cars	CUB-200-2011	89.50	
Ours	Stanford Cars	Stanford Cars	89.44	
Ours	Stanford Cars	CIFAR-100	88.90	
Ours	Stanford Cars	CIFAR-10	88.63	
Independent	CIFAR-100	-	73.16	
DML	CIFAR-100	-	73.61	
Ours	CIFAR-100	Stanford Dogs	72.19	
Ours	CIFAR-100	CUB-200-2011	73.66	
Ours	CIFAR-100	Stanford Cars	72.53	
Ours	CIFAR-100	CIFAR-100	74.18	
Ours	CIFAR-100	CIFAR-10	73.37	
Independent	CIFAR-10	-	93.99	
DML	CIFAR-10	-	93.97	
Ours	CIFAR-10	Stanford Dogs	93.87	
Ours	CIFAR-10	CUB-200-2011	94.09	
Ours	CIFAR-10	Stanford Cars	93.93	
Ours	CIFAR-10	CIFAR-100	94.37	
Ours	CIFAR-10	CIFAR-10	94.15	

Table 3: Ensemble accuracy of reused two-node graphs optimized on another dataset [%].

Method	Optimizing	Number of nodes			
Method	Graph	2	3	4	5
Independent	-	70.09	71.41	72.06	72.32
Ours	Stanford Dogs	73.86	73.41	74.16	74.14
Ours	CUB-200-2011	72.43	73.78	74.03	74.55
Ours	Stanford Cars	72.79	73.35	72.86	74.56

Table 4: Ensemble accuracy on Stanford Dogs [%].

Table 5: Ensemble accuracy on CUB-200-2011 [%].

Method	Optimizing	Number of nodes			
Method	Graph	2	3	4	5
Independent -		65.26	65.27	66.40	66.66
Ours	Stanford Dogs	72.06	71.82	73.03	72.13
Ours	CUB-200-2011	69.81	74.17	71.27	74.05
Ours	Stanford Cars	71.19	72.56	70.26	73.43

Table 6: Ensemble accuracy on Stanford Cars [%].

Method	Optimizing	Number of nodes			es
method	Graph	2	3	4	5
Independent	-	88.49	89.23	89.55	89.48
Ours	Stanford Dogs	89.76	89.94	89.98	90.41
Ours	CUB-200-2011	89.58	90.39	89.95	90.81
Ours	Stanford Cars	89.44	90.04	89.57	90.73



Fig. 6: Attention map of five nodes optimized by Stanford Dogs trained on other datasets with dimensionality reduction by UMAP.

8 N. Okamoto et al.

Used for ensemble					Fncomblo
node 1	node 2	node 3	node 4	node 5	Ensemble
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	74.50
\checkmark	\checkmark	\checkmark	\checkmark		74.30
\checkmark	\checkmark	\checkmark		\checkmark	74.06
\checkmark	\checkmark		\checkmark	\checkmark	74.31
\checkmark		✓	✓	\checkmark	74.23
	\checkmark	\checkmark	\checkmark	\checkmark	74.44
\checkmark	\checkmark	\checkmark			73.76
\checkmark	\checkmark		✓		73.80
\checkmark	\checkmark			 ✓ 	73.71
\checkmark		\checkmark	\checkmark		73.87
\checkmark		√		 ✓ 	73.28
\checkmark			✓	✓	73.65
	\checkmark	\checkmark	\checkmark		73.51
	\checkmark	✓		\checkmark	73.88
	\checkmark		✓	✓	73.40
		\checkmark	\checkmark	\checkmark	73.01
\checkmark	\checkmark				72.54
\checkmark		\checkmark			72.69
\checkmark			\checkmark		72.74
\checkmark				\checkmark	72.28
	\checkmark	✓			72.62
	\checkmark		✓		72.69
	\checkmark			\checkmark	72.18
		\checkmark	\checkmark		72.56
		\checkmark		\checkmark	71.52
			\checkmark	\checkmark	72.44

Table 7: Comparison of ensemble accuracy by nodes used in ensemble. [%].

6 Relationship between network nodes and ensemble accuracy

On the basis of the trend in Fig. 7 of the main paper, training with the five-node graph automatically designed using Stanford Dogs (Fig. 6d of the main paper) lead to a different attention map being acquired for each network node. As a result, node 2 became a network that strongly focused on the background, which may have a negative impact on the ensemble. In this section, we show the impact on the ensemble accuracy of each network node trained with the five-node graph automatically designed using Stanford Dogs. We used ABN based on ResNet as the network and Stanford Dogs as the dataset.

Table 7 shows the ensemble accuracy using selected network nodes. " \checkmark " means the network node used in the ensemble. The bold text in the Ensemble column shows the highest ensemble accuracy for each number of network nodes. The ensemble accuracy was highest when all network nodes were used. Therefore, we believe that the five-node graph is a training method that achieves high performance by acquiring different attention maps among network nodes that cooperate with each other.

References

- 1. Fukui, H., Hirakawa, T., Yamashita, T., Fujiyoshi, H.: Attention branch network: Learning of attention mechanism for visual explanation. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (2019)
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 770–778 (2016)
- 3. Khosla, A., Jayadevaprakash, N., Yao, B., Fei-Fei, L.: Novel dataset for fine-grained image categorization. In: First Workshop on Fine-Grained Visual Categorization, IEEE Conference on Computer Vision and Pattern Recognition (2011)
- Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3d object representations for fine-grained categorization. In: 4th International IEEE Workshop on 3D Representation and Recognition. Sydney, Australia (2013)
- 5. Krizhevsky, A., Hinton, G.: Learning multiple layers of features from tiny images. Tech. rep., Citeseer (2009)
- McInnes, L., Healy, J., Saul, N., Grossberger, L.: Umap: Uniform manifold approximation and projection. The Journal of Open Source Software 3(29), 861 (2018)
- Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011 dataset. Tech. Rep. CNS-TR-2011-001, California Institute of Technology (2011)
- 8. Zagoruyko, S., Komodakis, N.: Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In: International Conference on Learning Representations (2017)
- Zhang, Y., Xiang, T., Hospedales, T.M., Lu, H.: Deep mutual learning. In: IEEE Conference on Computer Vision and Pattern Recognition (2018)