



# Rethinking Ensemble-Distillation for Semantic Segmentation Based Unsupervised Domain Adaptation

Chen-Hao Chao, Bo-Wun Cheng, and Chun-Yi Lee

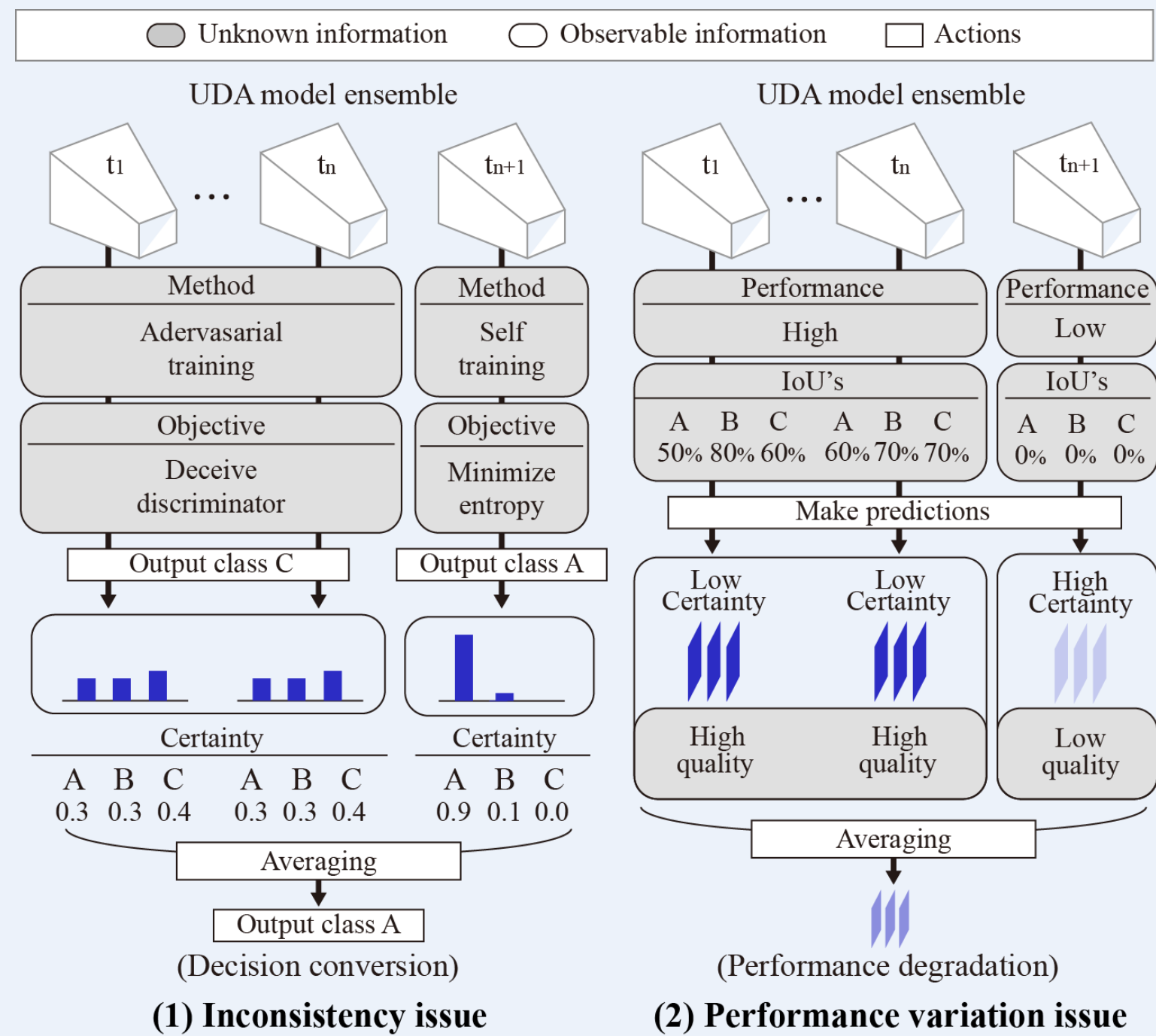
Elsa Lab, Department of Computer Science, National Tsing Hua University, Hsinchu, Taiwan



## Abstract

Recent researches on unsupervised domain adaptation (UDA) have demonstrated that end-to-end ensemble learning frameworks serve as a compelling option for UDA tasks. Nevertheless, these end-to-end ensemble learning methods often lack **flexibility** as any modification to the ensemble requires retraining of their frameworks. To address this problem, we propose a flexible ensemble-distillation framework for performing semantic segmentation based UDA, allowing any arbitrary composition of the members in the ensemble while still maintaining its superior performance. To achieve such flexibility, our framework is designed to be **robust** against the output inconsistency and the performance variation of the members within the ensemble. To examine the effectiveness and the robustness of our method, we perform an extensive set of experiments on both GTA5 to Cityscapes and SYNTHIA to Cityscapes benchmarks to quantitatively inspect the improvements achievable by our method. We further provide detailed analyses to validate that our design choices are practical and beneficial. The experimental evidence validates that the proposed method indeed offer **superior performance, robustness and flexibility** in semantic segmentation based UDA tasks against contemporary baseline methods.

## The Issues of Previous Methods



### (1) Inconsistency Issue

Since each model in the ensemble can be trained independently using different methods, **the scale** of the output certainty values may **not be consistent** across the ensemble. This may result in a situation that few members' decisions with high certainty values dominate the entire ensemble's output.

### (2) Performance Variation Issue

Since the **performance** (either per-class or average accuracy) of each teacher model in the ensemble may **vary substantially**, few under-performing members in the ensemble may cause the quality of the combined prediction to degrade.

## Methodology

### (1) A Comparison of the Proposed Method and Previous Methods

To address the **inconsistency** and **performance variation** issues, we introduce a new ensemble-distillation framework, and illustrate it in Fig.1(a). The main difference between the proposed method and the previous ones lies in two aspects: **output unification** and **fusion function**.

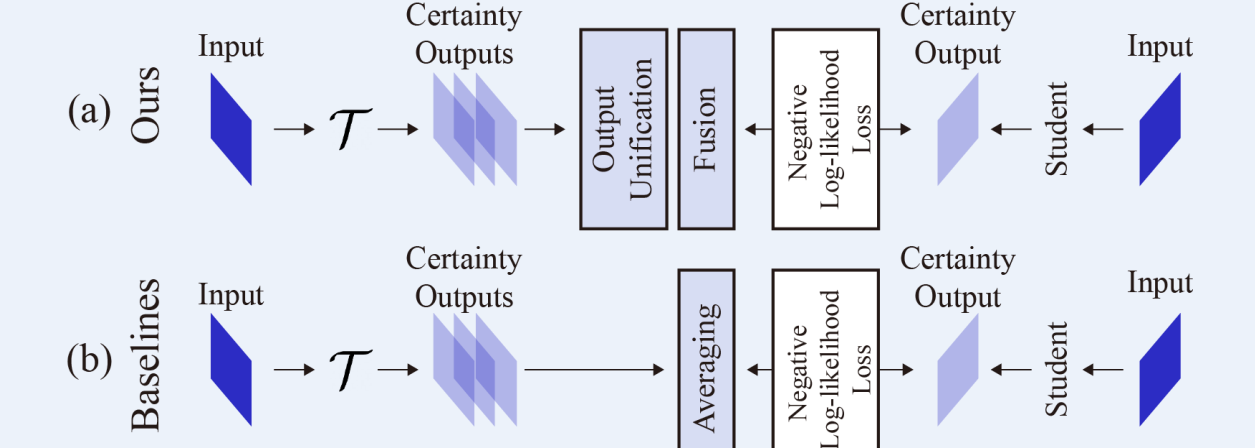


Figure 1: A comparison between (a) the proposed ensemble-distillation framework and (b) the baseline framework.

### (2) Output Unification

The soft predictions  $\hat{s}$  are unified first:

$$\hat{y}^{(p,c,t)} = \begin{cases} 1, & \text{if } c = \arg\max_{c \in \mathcal{C}} \{\hat{s}^{(p,c,t)}\} \\ 0, & \text{otherwise} \end{cases}$$

where  $\mathcal{C}$  is a set of semantic classes. The unification operation converts  $\hat{s}$  as **pseudo labels** to ensure that the subsequent fusion function can operate on items with a consistent scale.

### (3) Fusion Function

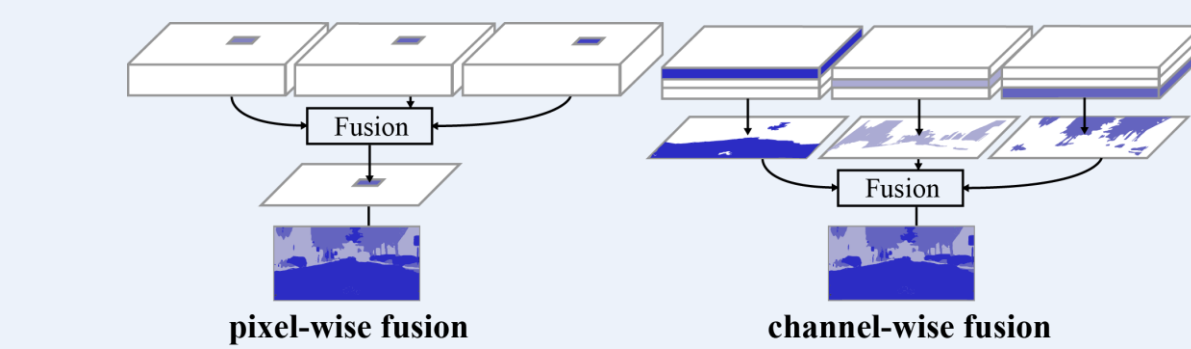


Figure 2: An illustration of the pixel-wise and channel-wise fusions.

The proposed **channel-wise fusion** rearranges the unified outputs according to a **fusion policy**  $\pi$  to construct the fused results. Fig.4. demonstrates the process of the rearrangement.

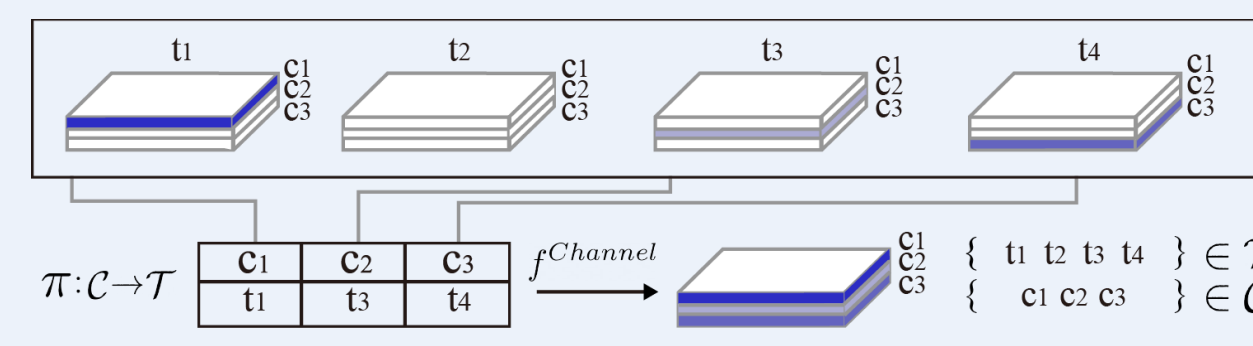


Figure 4: An illustrative example of  $\pi$  used in channel-wise fusion. Here,  $\mathcal{C}$  denotes the set of semantic classes and  $\mathcal{T}$  represents the set of teacher models in the ensemble.

### (4) Policy Selection Strategy



Figure 5: An illustration of how the quality of a teacher model's pseudo labels can affect the output certainty values of the student model.

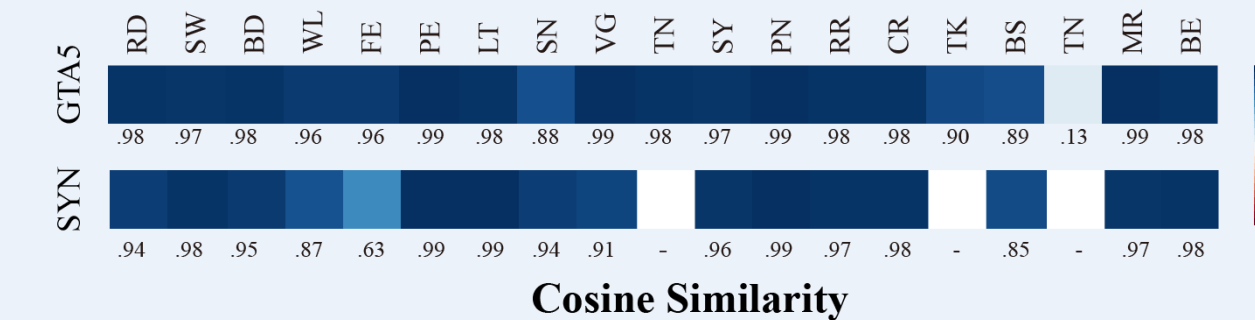


Figure 6: The cosine similarity between the per-class IoU of each  $t \in \mathcal{T}$  and the per-class certainty values of the student evaluated on the training set of Cityscapes.

To find an appropriate  $\pi$ , we adopt a certainty-aware policy selection strategy. The strategy is designed based on the observation that the student's output **certainty** values is highly **correlated** to the **quality** of the fused pseudo labels as shown in Fig.5. and Fig.6. Therefore, our framework first performs knowledge distillation on  $t \in \mathcal{T}$  and transfers their knowledge to  $|\mathcal{T}|$  identical student models using the unified outputs. Then, their output certainty values are measured to obtain the approximated performance of their corresponding teacher models. Finally, the policy that maximizes the student model's output certainty values is selected.

## Experimental Results

### (1) Performance on the Benchmarks

Method	Road	SideW	Build	Wall	Fence	Pole	Light	Sign	Veg	Termin	Sky	Person	Rider	Car	Truck	Bus	Train	Motor	Bike	mIoU
APDA [20]	85.6	52.8	79.0	29.5	25.5	26.8	34.6	19.9	83.7	40.6	77.9	59.2	28.3	84.6	34.6	49.2	8.0	32.6	39.6	45.9
PatchAlign [21]	92.3	51.9	82.1	29.2	25.1	24.5	33.8	33.0	82.4	32.8	82.2	58.6	27.2	84.3	33.4	46.3	2.2	29.5	32.3	46.5
AdvEnt [22]	89.4	53.1	81.0	26.6	26.8	27.2	33.5	24.7	83.9	36.7	78.8	58.7	30.5	84.8	38.5	44.5	1.7	31.6	32.4	45.5
IDA-MIT [26]	92.5	53.3	82.4	26.5	27.6	36.4	40.6	38.9	82.3	39.8	78.0	62.6	34.4	84.9	34.1	53.1	16.9	27.7	46.4	50.5
PT [35]	87.5	43.4	78.8	31.2	30.2	36.3	39.9	42.0	79.2	37.1	79.3	65.4	37.5	83.2	46.0	45.6	23.7	23.5	49.9	50.6
CBST [27]	91.8	53.5	80.5	32.7	21.0	34.0	28.9	20.4	83.9	34.2	80.9	53.1	24.0	82.7	30.3	35.9	16.0	25.9	42.8	45.9
MRKL [28]	91.0	55.4	80.0	33.7	21.4	37.3	32.9	24.5	85.0	34.1	80.8	57.7	24.6	84.1	27.8	30.1	26.9	26.0	42.3	47.1
R-MRNet [29]	90.4	51.2	85.1	36.9	25.6	37.5	48.8	48.5	85.3	34.8	81.1	64.4	36.8	86.3	34.9	52.2	1.7	29.0	44.6	50.3
DACS [30]	89.90	59.66	87.87	30.71	39.52	38.52	46.43	52.79	87.98	43.96	88.76	67.20	35.78	84.45	45.73	50.19	0.00	27.25	33.96	52.14
Source Only	37.40	21.43	56.80	4.93	22.14	32.38	34.62	24.90	78.96	11.92	63.71	55.55	13.83	88.11	21.99	29.78	2.36	28.41	33.96	34.80
End [43]	92.17	53.12	84.85	24.77	29.76	40.38	40.98	49.35	86.21	42.85	79.74	62.79	35.98	85.72	42.10	44.45	0.26	28.27	51.80	51.34
End* [50]	92.39	53.84	85.34	24.51	30.53	40.28	42.40	50.28	86.19	43.39	80.55	63.26	36.75	86.15	43.95	43.91	0.20	30.17	53.22	51.96
Ours (Pixel)	92.29	57.34	84.09	36.75	29.17	41.37	48.98	42.26	86.91	79.95	82.81	66.29	37.42	86.94	35.21	44.82	1.48	40.78	53.02	53.26
Ours (Channel)	94.43	66.90	88.07	39.46	41.80	43.24	49.08	56.00	88.01	45.83	87.79	67.58	38.05	90.08	57.64	51.90	0.00	46.57	55.28	57.98

Table 1: The quantitative results evaluated on the GTA5 to Cityscapes and SYNTHIA to Cityscapes benchmarks. The numbers presented in the middle and the last two columns correspond to per-class IoUs, mIoU, and mIoU\*, respectively.

The models used in our semantic segmentation based UDA model ensemble  $\mathcal{T}$  are highlighted in blue. The setting 'Source Only' indicates that the student model is trained only with the source domain ground truth annotations.

### (2) Robustness

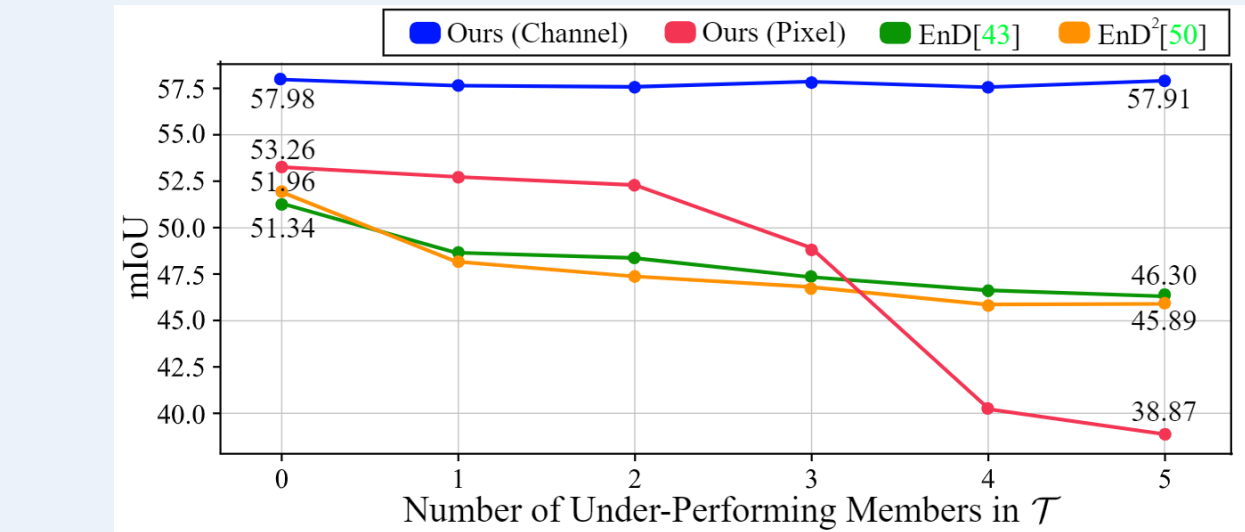


Table 2: The performance comparison of our framework with channel-wise fusion, our framework with pixel-wise fusion, and the baseline methods, with under-performing members ('Source Only' in Table 1) added to  $\mathcal{T}$ .

### (3) Flexibility

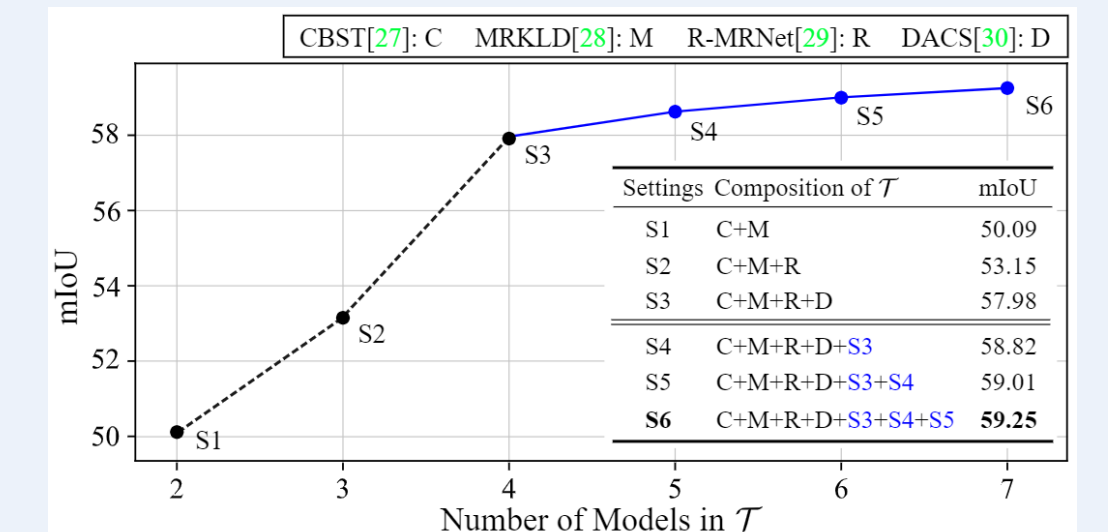


Table 3: The performance of the student models trained with the proposed framework using channel-wise fusion with different compositions of  $\mathcal{T}$ .

## Conclusion

In this paper, we presented a flexible ensemble-distillation framework to address the common pitfalls, i.e., the lack of robustness, of previous methods. We incorporated an output unification operation into the proposed framework to ensure that the fused outputs of the ensemble are free of the influence from the certainty inconsistency among the models in the ensemble. In addition, to tackle the performance variation issue, we proposed a channel-wise fusion function that is robust against this issue. As our framework is able to integrate different types of UDA methods while maintaining its robustness, it therefore pioneers a new direction for future semantic segmentation based UDA researches.

### Contributions

We introduce a flexible UDA ensemble-distillation framework which is robust against the inconsistency in the scale of the output certainty values and the performance variations among the members in an ensemble.

### Questions?

lance\_chao@gapp.nthu.edu.tw  
bobcheng15@gapp.nthu.edu.tw  
cylee@gapp.nthu.edu.tw

### Acknowledgements

TWCC NVIDIA  
MEDIATEK MOST 科技部  
Ministry of Science and Technology