# Appendix

## A Further Explanation of RAV

We further explain the RAV shown in Eq. (**??**) to elucidate the role of RAV in measuring the target sample value. The query function RAV can be expressed as:

$$
\begin{aligned}
\text{RAV}(x^{ut}) &= -\sum_{c=1}^{K}(p_c^{ut} + \widetilde{p}_c^{ut})\log(p_c^{ut}\widetilde{p}_c^{ut}) \\
&= -\sum_{c=1}^{K}(p_c^{ut} + \widetilde{p}_c^{ut})(\log p_c^{ut} + \log \widetilde{p}_c^{ut}) \\
&= -\sum_{c=1}^{K} p_c^{ut}\log p_c^{ut} - \sum_{c=1}^{K} p_c^{ut}\log \widetilde{p}_c^{ut} \quad (1) \\
&\quad -\sum_{c=1}^{K}\widetilde{p}_c^{ut}\log p_c^{ut} - \sum_{c=1}^{K}\widetilde{p}_c^{ut}\log \widetilde{p}_c^{ut} \\
&= \mathcal{H}(p^{ut}) + \mathcal{H}(\widetilde{p}^{ut}) \\
&\quad + \mathcal{H}(p^{ut}, \widetilde{p}^{ut}) + \mathcal{H}(\widetilde{p}^{ut}, p^{ut})
\end{aligned}
$$

where $\mathcal{H}(p^{ut}) + \mathcal{H}(\widetilde{p}^{ut})$ is used to measure whether the model can confidently predict the original and reconstructed features, and $\mathcal{H}(p^{ut}, \widetilde{p}^{ut}) + \mathcal{H}(\widetilde{p}^{ut}, p^{ut})$ is used to measure the discrepancy between the prediction probabilities of the original and reconstructed features. If $x^{ut}$ is not target domain-specific, that is, $x^{ut}$ has a high similarity with the source domain. Therefore, the model can make a high-confidence prediction for $x^{ut}$ and the semantic information of the original feature is the same as that contained in the reconstructed feature. Thus, the confidence of the prediction probability of $x^{ut}$ is high ($\mathcal{H}(p^{ut}) + \mathcal{H}(\widetilde{p}^{ut})$ is small), and the discrepancy between the prediction probabilities of the original feature and the reconstruction feature of $x^{ut}$ is small ($\mathcal{H}(p^{ut}, \widetilde{p}^{ut}) + \mathcal{H}(\widetilde{p}^{ut}, p^{ut})$ is small). Only when the model produces consistent and completely certain predictions for the original and reconstructed features (e.g. [1,0,0] and [1,0,0]), RAV reaches the minimum value of zero. For valuable target samples, namely those containing a large amount of target domain-specific knowledge, their RAVs are relatively high.

## B More Implementation Details of RASC-Ob

Previous ADA methods require experts to label the selected sample correctly, which requires $\log_2 K$ bits of information. In RASC-Ob, one-bit annotation is used, which only needs to determine whether the prediction of the selected sample is correct. Therefore, labeling a sample only requires one bit of information. For a fair comparison, we keep the amount of information obtained by one-bit annotation the same as that obtained by the traditional annotation method. Therefore, under the same labeling budget, one-bit annotation selects $\log_2 K$ times more samples than traditional annotation methods. That is, RASC-Ob selects $n_{ob} = \log_2 K \cdot B \cdot |\mathcal{D}_{ut}|$ target samples in total for annotation.

| Method | A→D | A→W | D→A | D→W | W→A | W→D | Avg |
|---|---|---|---|---|---|---|---|
| Source Only | 81.5 | 75.0 | 63.1 | 95.2 | 65.7 | 99.4 | 80.0 |
| Random | 90.9 | 90.4 | 80.4 | 98.3 | 80.8 | _99.6_ | 90.1 |
| AADA | 93.5 | 93.1 | 83.2 | 99.7 | 84.2 | **100.0** | 92.3 |
| TQS | 96.4 | 96.4 | 86.4 | **100.0** | 87.1 | **100.0** | _94.4_ |
| SDM-AG | **98.8** | 97.6 | **89.6** | **100.0** | **88.0** | **100.0** | **95.6** |
| TL-ADA | 97.8 | _97.9_ | 85.0 | 99.8 | 85.3 | **100.0** | 94.3 |
| **RASC** | 98.0 | **98.6** | _89.2_ | **100.0** | _87.9_ | **100.0** | **95.6** |
| **RASC-Ob** | _98.6_ | 96.4 | 85.4 | _99.9_ | 85.9 | **100.0** | _94.4_ |

Table 5: Comparison results (Accuracy: %) on Office-31 with 10% labeling budget.

## C Additional Comparative Results with Larger Labeling Budget

The main purpose of ADA is to maximize adaptation performance with the smallest possible labeling budget. That is, a good ADA method should demonstrate satisfactory performance under a small labeling budget. In the main paper, we demonstrate the superiority of our RASC and RASC-Ob under the small labeling budget (5%) through experiments. In addition, to further investigate the performance of our method under larger labeling budgets, we double the labeling budget in the main paper, conduct experiments under a 10% labeling budget, and compare with existing ADA methods.

**Office-31.** Results on Office-31 are shown in Table 5. Due to the small size of Office-31 and the minor discrepancy across domains, its performance approaches saturation when using a large labeling budget. Our RASC and RASC-Ob achieve optimal and suboptimal performance, respectively. It demonstrates that our sample selection strategy is suitable for small-scale datasets with large labeling budgets.

**Office-Home.** Results on Office-Home are shown in Table 6. We can observe that our RASC and RASC-Ob outperform other comparison methods significantly in all tasks. Especially on some difficult tasks, such as P→C and A→C. It demonstrates that our method can still select the most valuable target samples to continuously improve adaptation performance under a large labeling budget. Moreover, we can observe that RASC-Ob using one-bit annotation is slightly better than RASC using the traditional annotation method, demonstrating that introducing one-bit annotation into ADA is a promising attempt.

**VisDA.** We also conduct experiments on the large-scale dataset, and the results are shown in the first column of Table 6. Our RASC outperforms other methods, and the performance of RASC-Ob using a simpler annotation method is second only to SDM-AG. It demonstrates that our method is also suitable for large-scale datasets with large inter-domain discrepancy under larger labeling budgets.

| Method | VisDA | Office-Home | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Syn→Real | A→C | A→P | A→R | C→A | C→P | C→R | P→A | P→C | P→R | R→A | R→C | R→P | Avg |
| Source Only | 44.7 | 42.1 | 66.3 | 73.3 | 50.7 | 59.0 | 62.6 | 51.9 | 37.9 | 71.2 | 65.2 | 42.6 | 76.6 | 58.3 |
| Random | 82.1 | 61.6 | 80.5 | 80.1 | 61.8 | 76.8 | 73.9 | 64.1 | 60.2 | 78.4 | 72.9 | 62.6 | 84.8 | 71.5 |
| AADA | 84.6 | 65.8 | 84.5 | 82.2 | 64.1 | 80.6 | 76.1 | 67.6 | 62.6 | 80.1 | 73.7 | 66.1 | 88.6 | 74.3 |
| TQS | 87.2 | 68.0 | 87.7 | 85.7 | 67.0 | 83.0 | 78.7 | 69.3 | 64.5 | 83.9 | 77.8 | 68.9 | 90.6 | 77.1 |
| SDM-AG | 89.3 | 71.2 | 89.8 | 88.3 | 71.4 | 86.3 | 83.0 | 73.9 | 69.1 | 86.2 | 81.8 | 71.6 | 92.7 | 80.4 |
| TL-ADA | 90.2 | 70.7 | 87.9 | 86.9 | 74.3 | 87.4 | 85.4 | 74.5 | 69.2 | 87.4 | 81.4 | 70.2 | 90.4 | 80.5 |
| **RASC** | **90.2** | **78.7** | **91.5** | 89.2 | 78.2 | 90.7 | 86.9 | 78.2 | **78.2** | 89.4 | 83.8 | **79.7** | **94.0** | 84.9 |
| **RASC-Ob** | 87.7 | 76.4 | 90.8 | **90.4** | **79.0** | **90.9** | **88.1** | **80.0** | 74.6 | **90.7** | **85.1** | **79.7** | 93.8 | **85.0** |

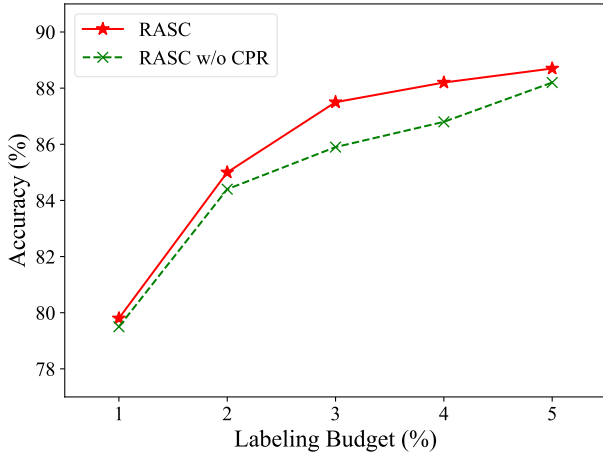Table 6: Comparison results (Accuracy: %) on VisDA and Office-Home with 10% labeling budget.



Figure 6: Performance with and without class prediction reliability-based sample rejection strategy (CPR) on VisDA.
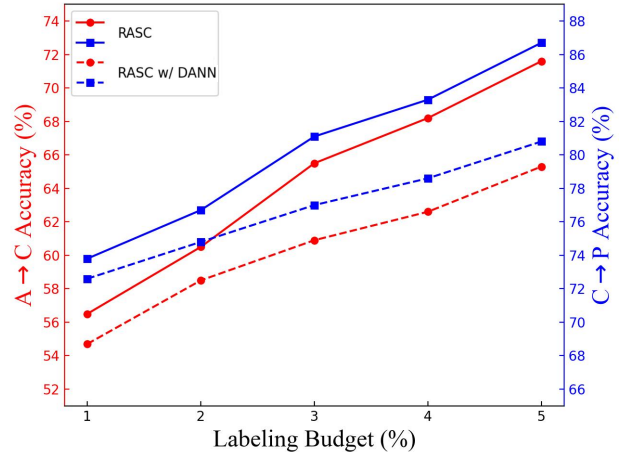


Figure 7: Comparison of different domain adaptation methods on Office-Home.

# D  Additional Analysis

## D.1  Effect of Class Prediction Reliability-Based Sample Rejection Strategy

In the main paper, we demonstrate that the class prediction reliability-based sample rejection strategy (CPR) not only considers class-level value but also helps to ensure sample diversity. To further investigate the effect of CPR on adaptation performance, we conduct experiments on VisDA. The results are shown in Figure 6, we can observe that the performance of RASC w/o CPR is always worse than that of RASC, which demonstrates the necessity of considering class-level value and ensuring sample diversity.

## D.2  Comparison of Domain Adaptation Methods

Unlike previous ADA methods that directly used the original DA methods, we design a more refined contrastive learning-based gradual active domain adaptation framework specifically for ADA. To further demonstrate its superiority, we compared it with the most commonly used DA method in ADA: DANN. We conduct experiments in A→C and C→P on Office-Home, and the results are shown in Figure 7. For a fair comparison, we only replace our contrastive learning-based gradual active domain adaptation framework with DANN for RASC w/ DANN, while keeping the remaining components of RASC unchanged. We can observe that RASC consistently outperforms RASC w/ DANN significantly. This is because our DA framework takes into account the unique domain differences in ADA. This demonstrates the superiority of our DA framework and the necessity of designing a DA framework specifically for ADA.
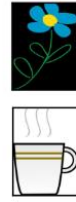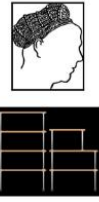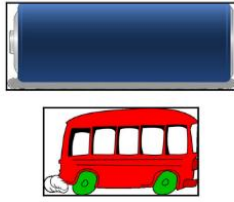
## D.3  Visualizing Selected Samples

To gain a more intuitive understanding of the characteristics of the selected target samples in RASC, we show the top 10 selected samples in RASC with the highest RAV in Figure 8. We can observe that the samples selected in each round are difficult to distinguish, and their styles differ greatly from the source domain, which helps the model learn more target domain-specific knowledge. In addition, although our method does not explicitly constrain the diversity of samples, it can be observed that the categories of samples selected in each round are different, which helps the model learn the discriminative features of each category in the target domain.
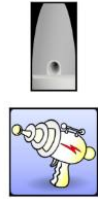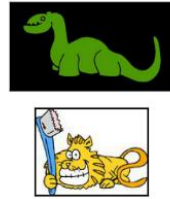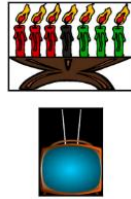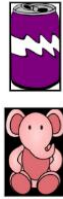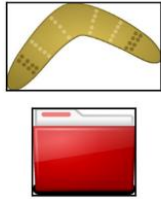
## D.4  T-SNE of Sample Selection Process in RASC

To further analyze the sample selection process of RASC, we visualize the target domain features after each round of sample selection by t-SNE. From Figure 9(a)-(e), we can observe
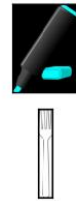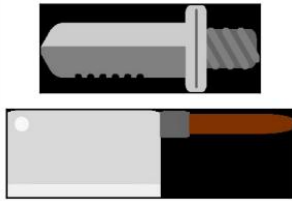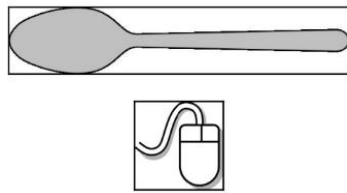
that the target samples selected in each round are located in the label-chaotic area, which means our reconfigurability-aware selection strategy can select hard-adaptive samples from a large number of target samples. Moreover, the samples selected in each round exhibit good diversity. In addition, we can observe that as the sample selection progresses, samples located in the label-chaotic area are gradually separated, ultimately forming good clustering in the target domain, as shown in Figure 9(f). It demonstrates that the samples selected in each round that are labeled and used for training can provide valuable information to the model to continuously improve the adaptation performance.
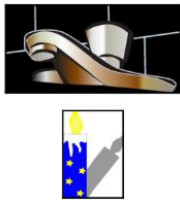
(a) Sampling Round 1
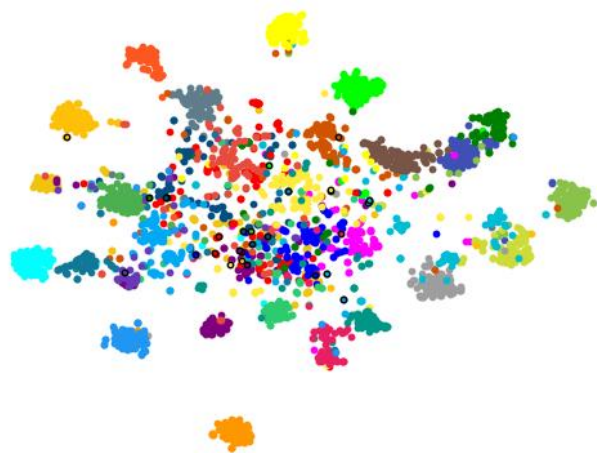

(b) Sampling Round 2


(c) Sampling Round 3


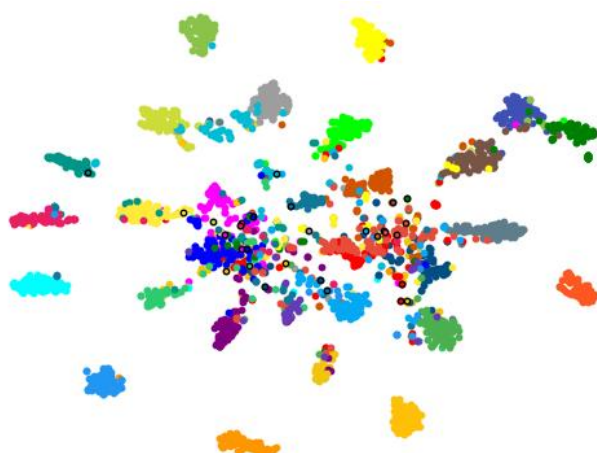(d) Sampling Round 4


(e) Sampling Round 5

Figure 8: The selected target samples in A→C on Office-Home. For each sampling round, we show the top 10 selected samples in RASC with the highest RAV.
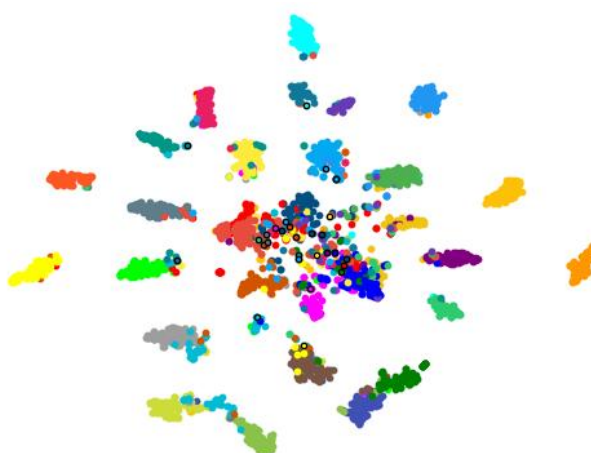
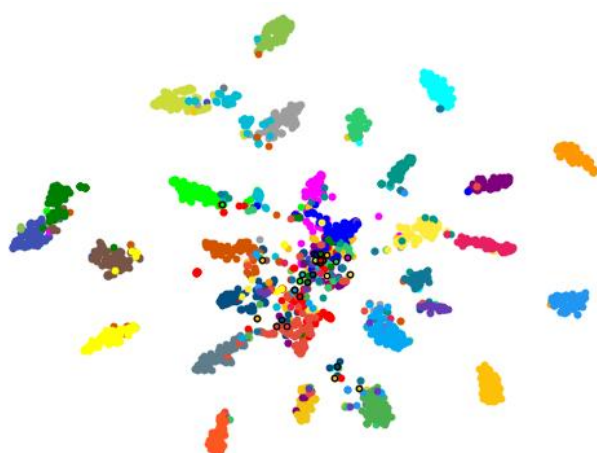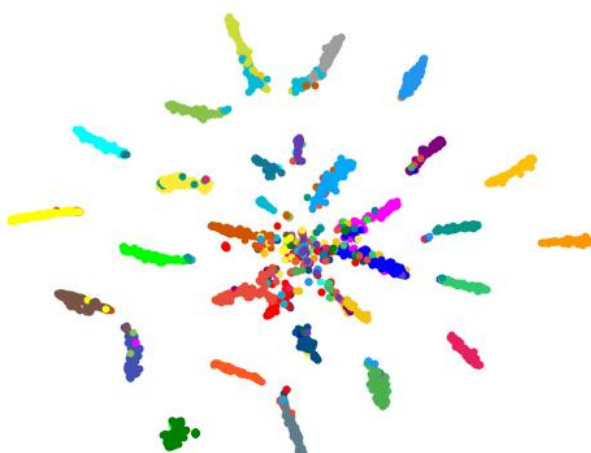Figure 9: T-SNE visualization of RASC sample selection process during training in W→A on Office-31. In the figure, each point is a target sample and the categories of samples are represented by different colors. During each round of sample selection, the selected target samples are surrounded with a black circle.