# Supplementary Material of "Online Learning for Noisy Labeled Streams"

Jinjie Qiu, Shengda Zhuo, Philip S.Yu,Changdong Wang, and Shuqiang Huang

**Outline.** This document supplements the manuscript in the following aspects. Section 1 introduces compared methods, and Section 2 presents the complete experimental results that supplement.

## 1 Compared Methods

This section presents complete descriptions of methods which used for comparison analysis.

**Compared Methods.** We compared OLNLS with 13 related state-of-the-art online learning (OL), offline noisy labels learning (OFNL), and online noisy labels learning (ONL$^2$).

- FOBOS (Singer and Duchi 2009) (OL): is an online learning algorithm that integrates gradient descent with regularization, designed to update models incrementally and promote sparsity by penalizing non-essential features.

- RDA (Singer and Duchi 2009) (OL): offers a base framework for direct training of online learners on observed features, utilizing a projected subgradient approach to enforce sparsity and eliminate features with negligible coefficients.

- FTRL (Singer and Duchi 2009) (OL): is an online learning algorithm that uses regularization techniques to accumulate all past gradients for parameter updates, which balances between aggressive learning and maintaining stability, especially in sparse data environments.

- FTRL-P (FTRL-Promixal) (Singer and Duchi 2009) (OL): is an online learning algorithm that extends FTRL by incorporating proximal regularization, improving convergence on sparse features and providing robustness against feature correlation, particularly in high-dimensional spaces.

- OVFM (He et al. 2021) (OL): manages mixed data streams by using Gaussian cointegration modeling to link observed and distributed data, enhancing model convergence despite mixed data challenges.

- OLI$^2$DS (You et al. 2023) (OL): optimizing feature selection via empirical risk minimization and using dynamic cost strategies for class imbalance to improve model efficacy.

- OSLMF (Wu et al. 2023) (OL): applying Gaussian Copula for imputation and Density Peak Clustering for label reconstruction, thus enriching the model with extensive label information.

- Co-teaching+ (Yu et al. 2019) (OFNL): addresses noisy label learning in offline settings, using a dual-model approach to pick samples with small loss and discrepant predictions for mutual updates, focusing on reliable samples to mitigate label noise impact.

- PHuber (Menon et al. 2019) (OFNL): employs robustness regularization in offline label noise learning to prevent overfitting to noise by modifying gradients, minimizing noisy label interference in the model.

- Bey.Ima. (BeyondImages) (Zhu, Wang, and Liu 2022) (OFNL): tackles label noise in offline learning by correcting loss through an information-theoretic method that extracts key features from low-dimensional representations and estimates a noise transition matrix for model compensation, diminishing noisy label interference.

- AdaStream (Zhang et al. 2022) (ONL$^2$): innovates in data stream mining under conditions of incomplete and noisy labels, using a landmark assumption to estimate the noise transition matrix and leveraging model reuse ensembles to bolster method stability.

- PA (Crammer et al. 2006) (ONL$^2$): is an online learning algorithm using first-order linear learning with a threshold for sample training to lessen sensitivity to noisy labels and prevent overfitting to noise distributions.

## 2 Complete Experiments for noise-label data streams

This section presents additional experimental results that were precluded of the main paper because of page limitation. These results include:

1. The complete CAR results yielded from all 13 datasets, which are shown in Table 1, Table 2, and Table 3. The results are collected from experiments in the setting of symmetric noisy labeled, asymmetric noise labeled, and flip noisy labeled data streams, compared with online learning algorithms.

2. The complete CAR trends yielded from all 13 datasets on symmetric noisy labeled, asymmetric noise labeled, and flip noisy labeled data streams, shown in Figure 1, Figure 2, and Figure 3. The results are collected from experiments in the setting of symmetric noisy labeled, asymmetric noise labeled, and flip noisy labeled data streams, compared with online learning algorithms.

3. The complete CAR results yielded from all 13 datasets, which are shown in Table 4, Table 5, and Table 6. The results are collected from experiments in the setting of symmetric noisy labeled, asymmetric noise labeled, and flip noisy labeled data streams, compared with noisy labels learning algorithms.

4. The complete CAR trends yielded from all 13 datasets on symmetric noisy labeled, asymmetric noise labeled, and flip noisy labeled data streams, shown in Figure 4, Figure 5, and Figure 6. The results are collected from experiments in the setting of symmetric noisy labeled, asymmetric noise labeled, and flipped noisy labeled data streams, compared with noisy labels learning algorithms.

# References

Crammer, K.; Dekel, O.; Keshet, J.; Shalev-Shwartz, S.; and Singer, Y. 2006. Online passive aggressive algorithms.

He, Y.; Dong, J.; Hou, B.-J.; Wang, Y.; and Wang, F. 2021. Online Learning in Variable Feature Spaces with Mixed Data. In *ICDM*, 181–190. IEEE.

Menon, A. K.; Rawat, A. S.; Reddi, S. J.; and Kumar, S. 2019. Can gradient clipping mitigate label noise? In *ICLR*.

Singer, Y.; and Duchi, J. C. 2009. Efficient learning using forward-backward splitting. *In NeurIPS*, 22.

Wu, D.; Zhuo, S.; Wang, Y.; Chen, Z.; and He, Y. 2023. Online semi-supervised learning with mix-typed streaming features. In *AAAI*, volume 37, 4720–4728.

You, D.; Xiao, J.; Wang, Y.; Yan, H.; Wu, D.; Chen, Z.; Shen, L.; and Wu, X. 2023. Online Learning From Incomplete and Imbalanced Data Streams. *IEEE Transactions on Knowledge and Data Engineering*.

Yu, X.; Han, B.; Yao, J.; Niu, G.; Tsang, I.; and Sugiyama, M. 2019. How does disagreement help generalization against label corruption? In *ICML*, 7164–7173. PMLR.

Zhang, Z.-Y.; Qian, Y.-Y.; Zhang, Y.-J.; Jiang, Y.; and Zhou, Z.-H. 2022. Adaptive Learning for Weakly Labeled Streams. In *KDD*, 2556–2564.

Zhu, Z.; Wang, J.; and Liu, Y. 2022. Beyond images: Label noise transition matrix estimation for tasks with lower-quality features. In *ICML*, 27633–27653. PMLR.

Table 1: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR $\pm$ Standard Variance) for 13 datasets in the case of Symmetric Noisy labeled data. • indicates the cases that our method loses the comparison. * showing the total number of wins and losses for OLNLS. OOT: Out Of Time (24 hours).
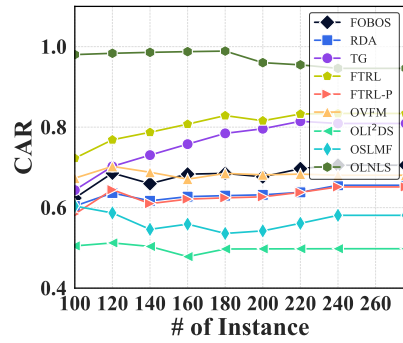
| Datasets | FOBOS | RDA | TG | FTRL | FTRL-P | OSLMF | OVFM | OLI$^2$DS | OLNLS |
|---|---|---|---|---|---|---|---|---|---|
| | \multicolumn{9}{c}{Symetric Noise $\rho_{-1} = \rho_{+1} = 0.4$} | | | | | | | | |
| breast | $.705 \pm .023$ | $.634 \pm .023$ | $.801 \pm .036$ | $.803 \pm .027$ | $.638 \pm .011$ | $.596 \pm .015$ | $.653 \pm .041$ | $.522 \pm .043$ | $\mathbf{.892 \pm .026}$ |
| dermatology | $.531 \pm .016$ | $.559 \pm .014$ | $.547 \pm .017$ | $.684 \pm .025$ | $.623 \pm .026$ | $.604 \pm .054$ | $.603 \pm .012$ | $.649 \pm .020$ | $\mathbf{.721 \pm .010}$ |
| wdbc | $.537 \pm .012$ | $.542 \pm .013$ | $.549 \pm .015$ | $.573 \pm .033$ | $.607 \pm .029$ | $.598 \pm .031$ | $.579 \pm .029$ | $.544 \pm .009$ | $\mathbf{.671 \pm .023}$ |
| diabetes | $.528 \pm .013$ | $.532 \pm .014$ | $.530 \pm .022$ | $.528 \pm .025$ | $.533 \pm .019$ | $.580 \pm .037$ | $.571 \pm .023$ | $.520 \pm .006$ | $\mathbf{.598 \pm .022}$ |
| german | $.508 \pm .017$ | $.515 \pm .013$ | $.513 \pm .009$ | $.570 \pm .017$ | $.534 \pm .015$ | $.578 \pm .017$ | $.607 \pm .022$ | $.534 \pm .024$ | $\mathbf{.694 \pm .006}$ |
| contraceptive | $.600 \pm .002$ | $.596 \pm .011$ | $.597 \pm .008$ | $.530 \pm .017$ | $.601 \pm .003$ | $.609 \pm .033$ | $.613 \pm .020$ | $.545 \pm .016$ | $\mathbf{.787 \pm .023}$ |
| splice | $.603 \pm .002$ | $.610 \pm .006$ | $.604 \pm .005$ | $.694 \pm .009$ | $.599 \pm .003$ | $.500 \pm .030$ | $.576 \pm .035$ | $.554 \pm .008$ | $\mathbf{.778 \pm .001}$ |
| mushroom | $.577 \pm .005$ | $.653 \pm .037$ | $.621 \pm .006$ | $.794 \pm .017$ | $.618 \pm .004$ | $.860 \pm .026\bullet$ | $.763 \pm .011$ | $.674 \pm .018$ | $.833 \pm .020$ |
| marketing | $.580 \pm .003$ | $.608 \pm .011$ | $.613 \pm .009$ | $.767 \pm .008$ | $.592 \pm .003$ | $.717 \pm .038$ | $.696 \pm .006$ | $.565 \pm .010$ | $\mathbf{.818 \pm .005}$ |
| hapt | $.593 \pm .001$ | $.687 \pm .009$ | $.623 \pm .006$ | $.671 \pm .007$ | $.600 \pm .009$ | $.575 \pm .019$ | $.596 \pm .005$ | $.555 \pm .004$ | $\mathbf{.773 \pm .008}$ |
| ring | $.526 \pm .008$ | $.529 \pm .008$ | $.525 \pm .008$ | $.530 \pm .006$ | $.528 \pm .008$ | $.599 \pm .012$ | $.566 \pm .016$ | $.525 \pm .006$ | $\mathbf{.656 \pm .005}$ |
| magic04 | $.598 \pm .001$ | $.595 \pm .001$ | $.599 \pm .002$ | $.718 \pm .005$ | $.599 \pm .003$ | $.777 \pm .071$ | $.837 \pm .013$ | $.553 \pm .002$ | $\mathbf{.859 \pm .015}$ |
| a8a | $.490 \pm .003$ | $.671 \pm .008$ | $.568 \pm .006$ | $.671 \pm .004$ | $.576 \pm .003$ | $.664 \pm .006$ | $.617 \pm .005$ | $.543 \pm .005$ | $\mathbf{.696 \pm .019}$ |
| loss/win | 0/13 | 0/13 | 0/13 | 0/13 | 0/13 | 1/12 | 0/13 | 0/13 | $\mathbf{1/103^*}$ |
| $p$-value | .0001 | .0001 | .0001 | .0001 | .0001 | .0004 | .0001 | .0001 | $--$ |
| $F$-rank | 2.6538 | 4.3462 | 4.0000 | 6.1538 | 4.6923 | 5.8462 | 5.9231 | 2.4615 | 8.9231 |

Table 2: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR $\pm$ Standard Variance) for 13 datasets in the case of ASymmetric Noisy labeled data. • indicates the cases that our method loses the comparison. * showing the total number of wins and losses for OLNLS. OOT: Out Of Time (24 hours).
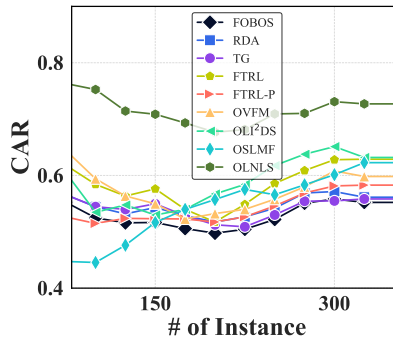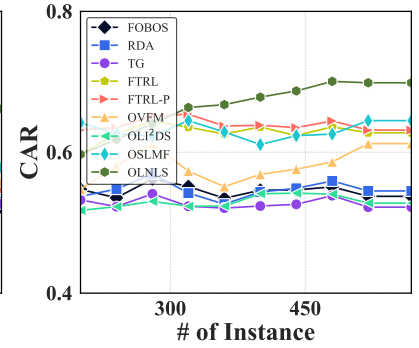
| Datasets | FOBOS | RDA | TG | FTRL | FTRL-P | OSLMF | OVFM | OLI$^2$DS | OLNLS |
|---|---|---|---|---|---|---|---|---|---|
| | \multicolumn{9}{c}{ASymmetric Noise $\rho_{-1} = 0.3, \rho_{+1} = 0.5$} | | | | | | | | |
| breast | $.752 \pm .019$ | $.662 \pm .007$ | $.800 \pm .017$ | $.781 \pm .014$ | $.678 \pm .018$ | $.591 \pm .037$ | $.656 \pm .013$ | $.533 \pm .033$ | $\mathbf{.812 \pm .020}$ |
| dermatology | $.588 \pm .010$ | $.598 \pm .008$ | $.588 \pm .015$ | $.717 \pm .031$ | $.659 \pm .031$ | $.610 \pm .027$ | $.673 \pm .023$ | $.657 \pm .024$ | $\mathbf{.749 \pm .018}$ |
| wdbc | $.569 \pm .005$ | $.568 \pm .012$ | $.578 \pm .008$ | $.602 \pm .038$ | $.589 \pm .031$ | $.588 \pm .047$ | $.616 \pm .016$ | $.556 \pm .011$ | $\mathbf{.677 \pm .027}$ |
| diabetes | $.486 \pm .014$ | $.484 \pm .016$ | $.493 \pm .009$ | $.483 \pm .015$ | $.493 \pm .011$ | $.510 \pm .035$ | $.492 \pm .016$ | $.510 \pm .015\bullet$ | $.510 \pm .021$ |
| german | $.491 \pm .018$ | $.495 \pm .016$ | $.497 \pm .014$ | $.470 \pm .007$ | $.503 \pm .013$ | $.476 \pm .033$ | $.487 \pm .039$ | $.536 \pm .019$ | $\mathbf{.556 \pm .074}$ |
| contraceptive | $.613 \pm .009$ | $.618 \pm .010$ | $.618 \pm .011$ | $.579 \pm .011$ | $.614 \pm .005$ | $.622 \pm .047$ | $.645 \pm .038$ | $.569 \pm .017$ | $\mathbf{.730 \pm .044}$ |
| splice | $.655 \pm .002$ | $.671 \pm .005$ | $.660 \pm .005$ | $.684 \pm .009$ | $.650 \pm .002$ | $.607 \pm .025$ | $.639 \pm .021$ | $.547 \pm .010$ | $\mathbf{.759 \pm .000}$ |
| mushroom | $.609 \pm .004$ | $.644 \pm .014$ | $.643 \pm .005$ | $.732 \pm .013$ | $.632 \pm .004$ | $.736 \pm .071$ | $.743 \pm .006$ | $.679 \pm .016$ | $\mathbf{.771 \pm .008}$ |
| marketing | $.683 \pm .003$ | $.745 \pm .002$ | $.720 \pm .002$ | $.820 \pm .001$ | $.673 \pm .003$ | $.767 \pm .017$ | $.783 \pm .006$ | $.574 \pm .017$ | $\mathbf{.826 \pm .002}$ |
| hapt | $.655 \pm .002$ | $.758 \pm .007$ | $.696 \pm .003$ | $.763 \pm .004$ | $.666 \pm .004$ | $.653 \pm .009$ | $.653 \pm .017$ | $.561 \pm .005$ | $\mathbf{.818 \pm .007}$ |
| ring | $.504 \pm .005$ | $.500 \pm .003$ | $.502 \pm .005$ | $.505 \pm .006$ | $.503 \pm .006$ | $.502 \pm .012$ | $.514 \pm .006$ | $.528 \pm .003$ | $\mathbf{.547 \pm .018}$ |
| magic04 | $.628 \pm .002$ | $.626 \pm .002$ | $.627 \pm .001$ | $.732 \pm .002$ | $.630 \pm .003$ | $.631 \pm .066$ | $.796 \pm .007$ | $.566 \pm .004$ | $\mathbf{.811 \pm .022}$ |
| a8a | $.460 \pm .003$ | $.671 \pm .013$ | $.575 \pm .007$ | $.740 \pm .012$ | $.638 \pm .005$ | $.740 \pm .005$ | $.684 \pm .005$ | $.547 \pm .004$ | $\mathbf{.757 \pm .000}$ |
| *loss/win* | 0/13 | 0/13 | 0/13 | 0/13 | 0/13 | 0/13 | 0/13 | 1/12 | $\mathbf{1/103^*}$ |
| $p$-value | .0001 | .0001 | .0001 | .0001 | .0001 | .0013 | .0001 | .0013 | $--$ |
| $F$-rank | 3.3462 | 4.0385 | 4.4615 | 5.8846 | 4.5769 | 4.7308 | 5.7308 | 3.3077 | 8.9231 |

Table 3: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR $\pm$ Standard Variance) for 13 datasets in the case of Flip Noisy labeled data. ● indicates the cases that our method loses the comparison. * showing the total number of wins and losses for OLNLS. OOT: Out Of Time (24 hours).

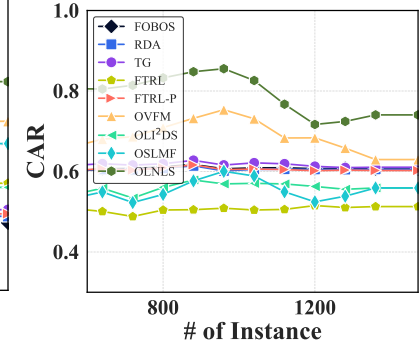| Datasets | | FOBOS | RDA | TG | FTRL | FTRL-P | OSLMF | OVFM | OLI$^2$DS | OLNLS |
|---|---|---|---|---|---|---|---|---|---|---|
| | | \multicolumn{9}{c}{Filp Noise Rate = 0.4} | | | | | | | | |
| breast | | $.677 \pm .020$ | $.617 \pm .032$ | $.736 \pm .055$ | $.790 \pm .037$ | $.641 \pm .021$ | $.582 \pm .070$ | $.656 \pm .013$ | $.561 \pm .056$ | $\mathbf{.896 \pm .052}$ |
| dermatology | | $.537 \pm .029$ | $.554 \pm .026$ | $.551 \pm .033$ | $.697 \pm .032$ | $.616 \pm .024$ | $.534 \pm .053$ | $.601 \pm .013$ | $.640 \pm .032$ | $\mathbf{.745 \pm .019}$ |
| wdbc | | $.541 \pm .009$ | $.539 \pm .007$ | $.537 \pm .010$ | $.564 \pm .016$ | $.579 \pm .026$ | $.574 \pm .018$ | $.576 \pm .020$ | $.543 \pm .012$ | $\mathbf{.673 \pm .023}$ |
| diabetes | | $.528 \pm .010$ | $.529 \pm .009$ | $.522 \pm .015$ | $.516 \pm .007$ | $.516 \pm .012$ | $.573 \pm .053$ | $.577 \pm .012$ | $.523 \pm .004$ | $\mathbf{.579 \pm .018}$ |
| german | | $.521 \pm .016$ | $.526 \pm .016$ | $.520 \pm .022$ | $.561 \pm .018$ | $.524 \pm .010$ | $.557 \pm .039$ | $.570 \pm .028$ | $.521 \pm .016$ | $\mathbf{.677 \pm .024}$ |
| contraceptive | | $.597 \pm .004$ | $.606 \pm .008$ | $.600 \pm .009$ | $.788 \pm .025$ | $.602 \pm .004$ | $.560 \pm .044$ | $.637 \pm .028$ | $.562 \pm .020$ | $\mathbf{.790 \pm .028}$ |
| splice | | $.600 \pm .001$ | $.605 \pm .004$ | $.604 \pm .004$ | $.690 \pm .004$ | $.596 \pm .003$ | $.526 \pm .031$ | $.573 \pm .027$ | $.551 \pm .009$ | $\mathbf{.779 \pm .001}$ |
| mushroom | | $.582 \pm .007$ | $.669 \pm .017$ | $.623 \pm .010$ | $.799 \pm .006$ | $.620 \pm .008$ | $.906 \pm .009$● | $.765 \pm .014$ | $.677 \pm .015$ | $.827 \pm .006$ |
| marketing | | $.578 \pm .007$ | $.626 \pm .022$ | $.619 \pm .005$ | $.771 \pm .009$ | $.592 \pm .006$ | $.736 \pm .019$ | $.694 \pm .012$ | $.567 \pm .008$ | $\mathbf{.815 \pm .006}$ |
| hapt | | $.592 \pm .002$ | $.700 \pm .017$ | $.622 \pm .004$ | $.673 \pm .004$ | $.599 \pm .005$ | $.568 \pm .011$ | $.587 \pm .012$ | $.552 \pm .003$ | $\mathbf{.772 \pm .014}$ |
| ring | | $.528 \pm .004$ | $.531 \pm .004$ | $.532 \pm .004$ | $.531 \pm .004$ | $.530 \pm .005$ | $.605 \pm .008$ | $.562 \pm .005$ | $.526 \pm .005$ | $\mathbf{.660 \pm .002}$ |
| magic04 | | $.599 \pm .001$ | $.597 \pm .001$ | $.602 \pm .003$ | $.719 \pm .002$ | $.600 \pm .002$ | $.745 \pm .039$ | $.830 \pm .007$ | $.554 \pm .003$ | $\mathbf{.869 \pm .014}$ |
| a8a | | $.489 \pm .003$ | $.670 \pm .017$ | $.569 \pm .003$ | $.681 \pm .008$ | $.577 \pm .001$ | $.676 \pm .008$ | $.614 \pm .004$ | $.544 \pm .003$ | $\mathbf{.711 \pm .002}$ |
| *loss/win* | | 0/13 | 0/13 | 0/13 | 0/13 | 0/13 | 1/12 | 0/13 | 0/13 | $\mathbf{1/103}^{*}$ |
| *p*-value | | .0001 | .0001 | .0001 | .0001 | .0001 | .0009 | .0001 | .0001 | $--$ |
| *F*-rank | | 3.0385 | 4.8077 | 4.0000 | 6.6154 | 4.1154 | 4.9231 | 6.0000 | 2.5769 | 8.9231 |



(a) breast



(b) dermatology



(c) wdbc



(d) diabetes



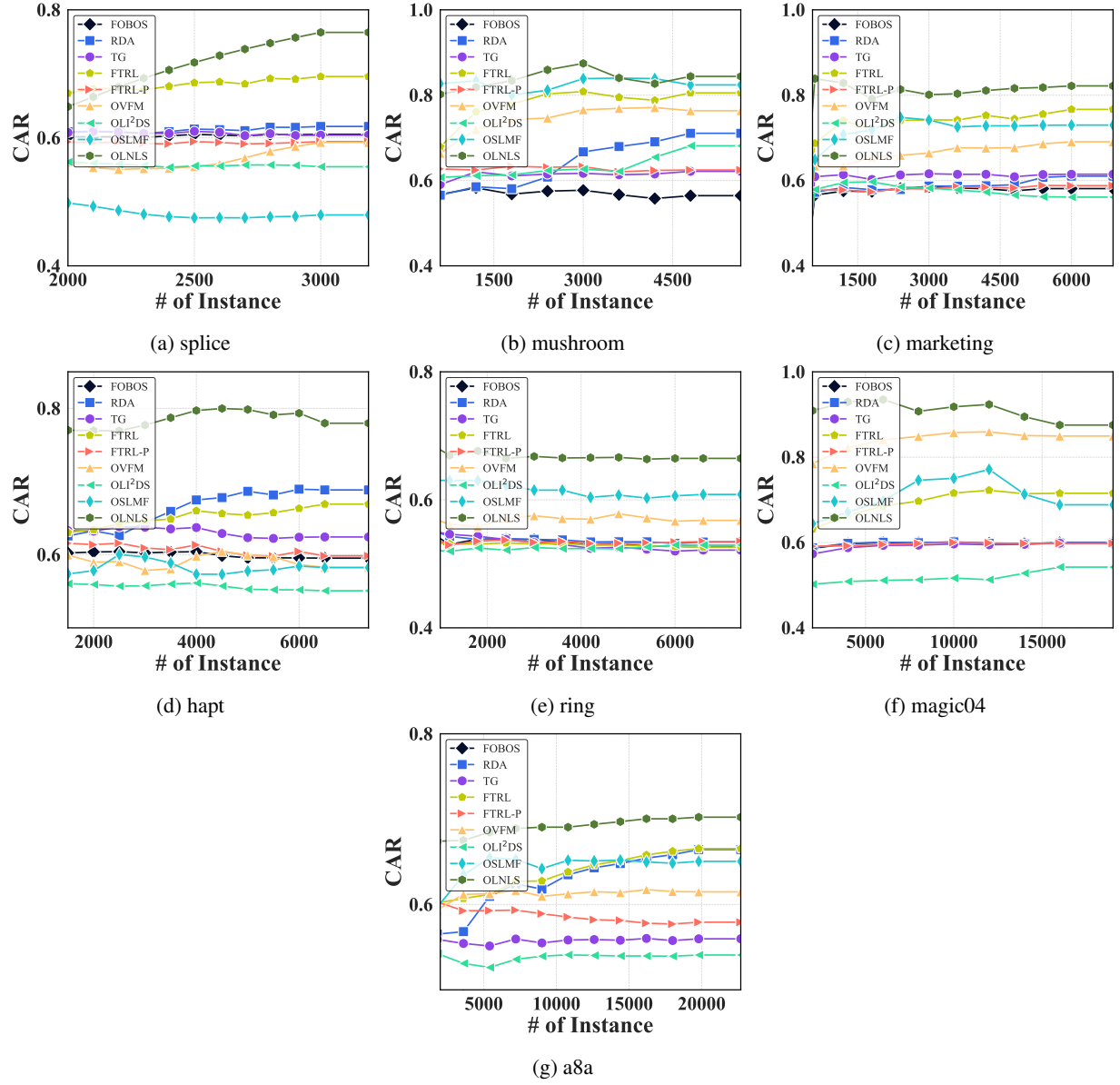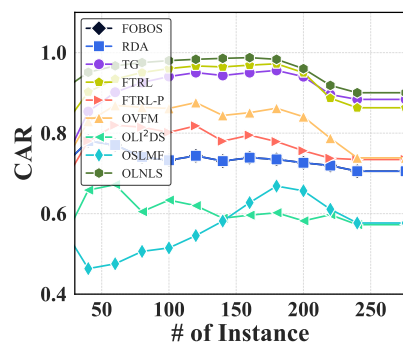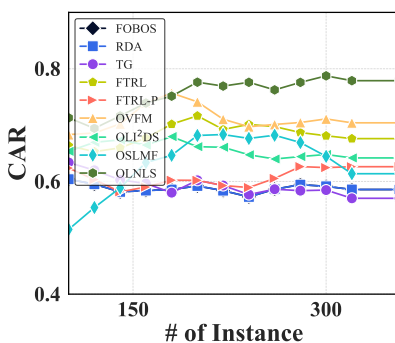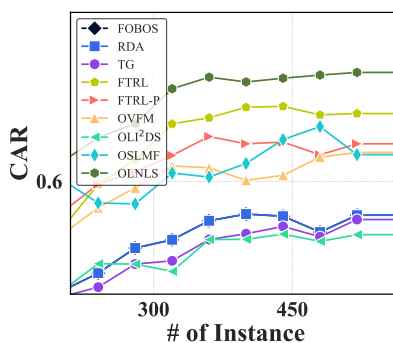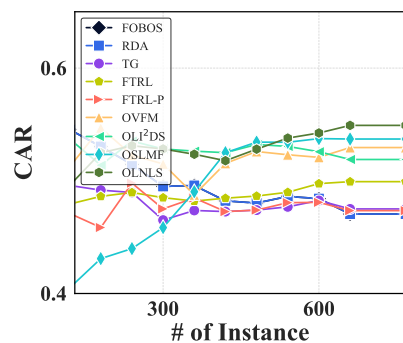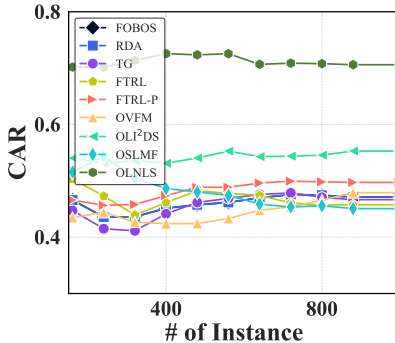(e) german



(f) contraceptive

Figure 1: The cumulative accuracy rate (CAR) trends of FOBOS, RDA, TG, FTRL, FTRL-P, OSLMF, OVFM, OL$^2$DS, and OLNLS in all 13 symmetric noisy labeled data streams.

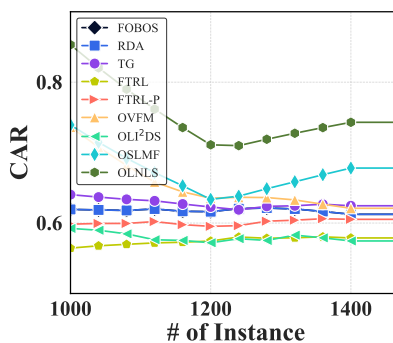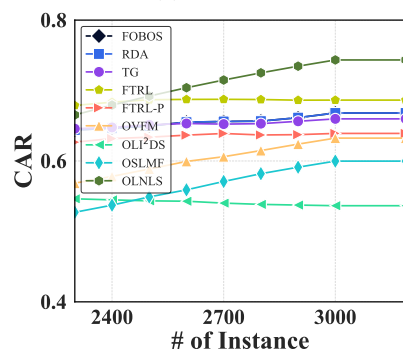(a) breast      (b) dermatology      (c) wdbc
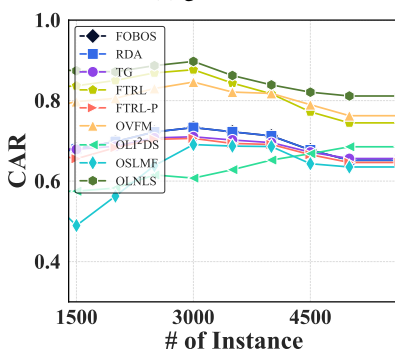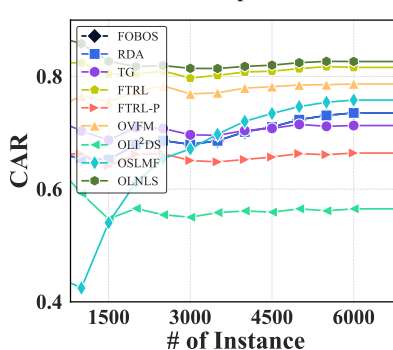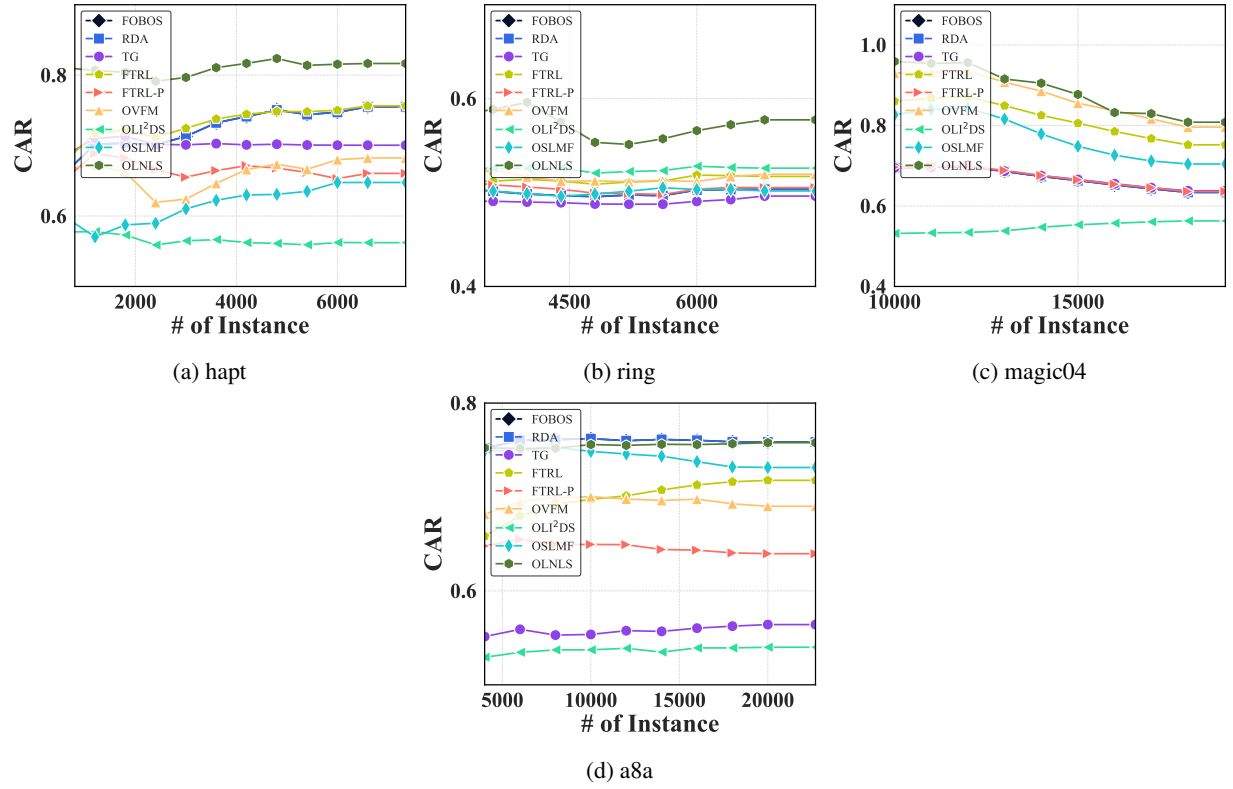
(d) diabetes      (e) german      (f) contraceptive
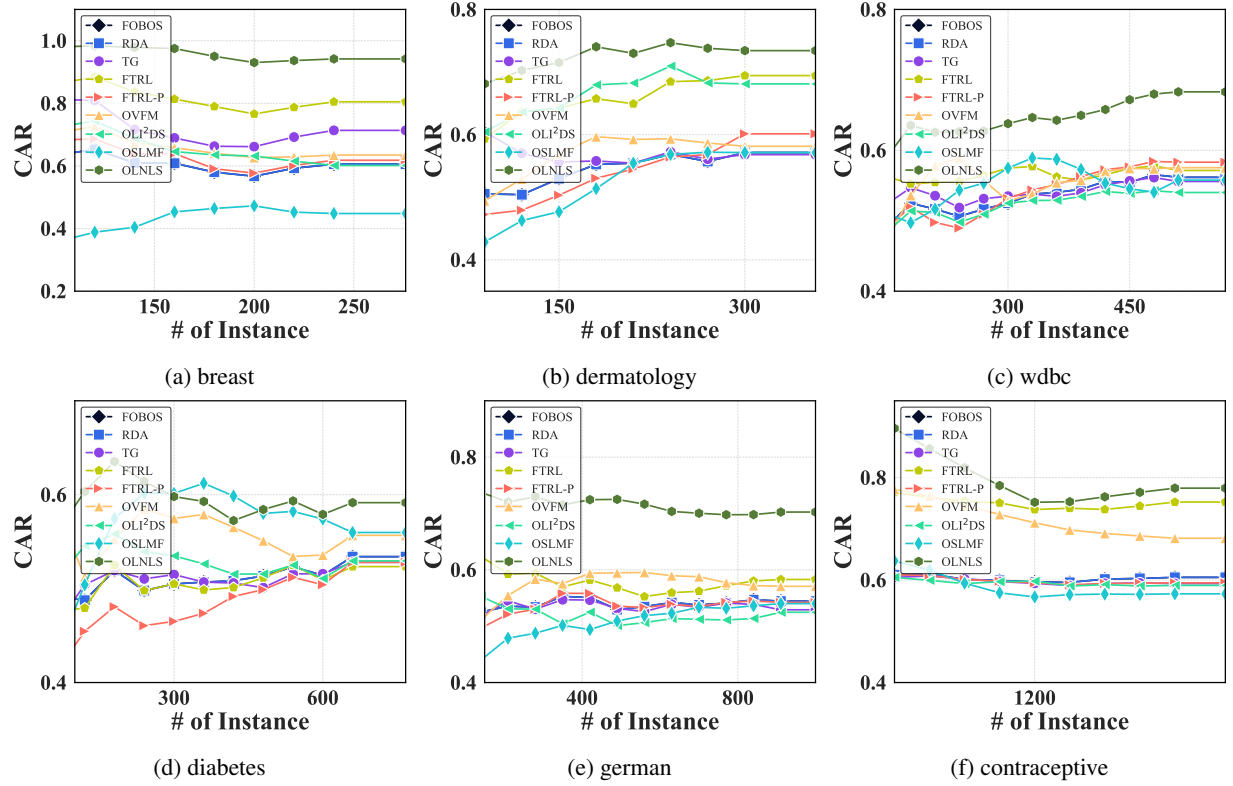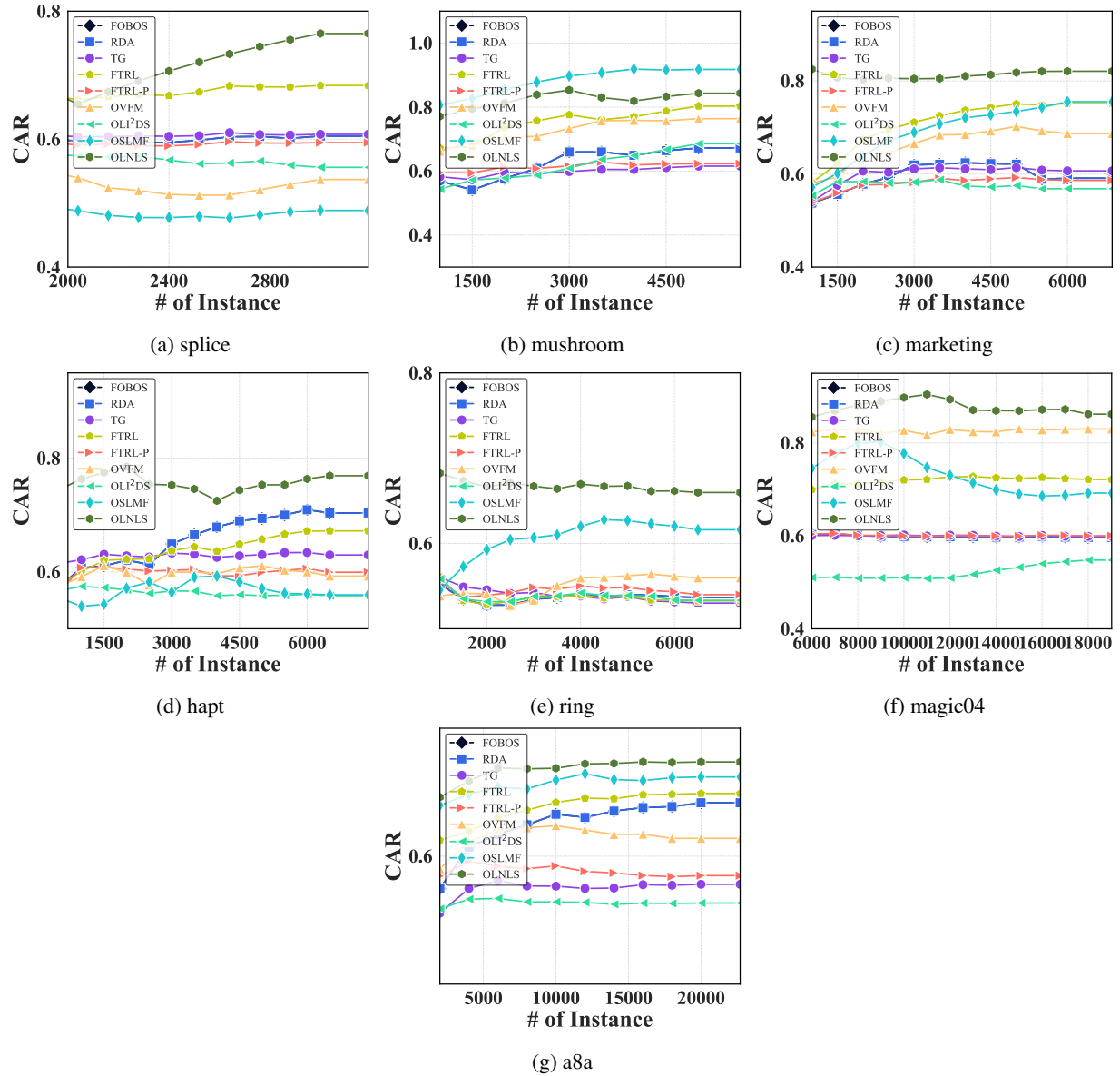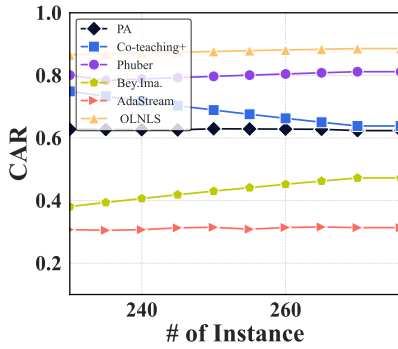
(g) splice      (h) mushroom      (i) marketing

Figure 2: The cumulative accuracy rate (CAR) trends of FOBOS, RDA, TG, FTRL, FTRL-P, OSLMF, OVFM, OL$^2$DS, and OLNLS in all 13 asymmetric noisy labeled data streams.

(a) breast

(b) dermatology

(c) wdbc

(d) diabetes

(e) german

(f) contraceptive

Figure 3: The cumulative accuracy rate (CAR) trends of FOBOS, RDA, TG, FTRL, FTRL-P, OSLMF, OVFM, $OL^2DS$, and OLNLS in all 13 flipped noisy labeled data streams.

Table 4: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR $\pm$ Standard Variance) for 13 datasets in the case of Symmetric Noisy labeled data. $\bullet$ indicates the cases that our method loses the comparison. * showing the total number of wins and losses for OLNLS. OOT: Out Of Time (24 hours).
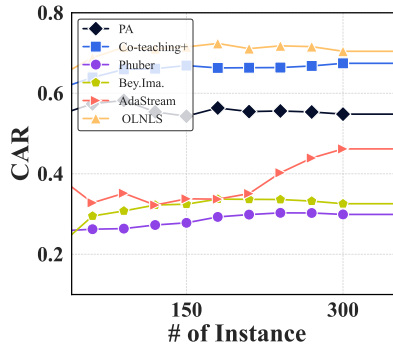
| | Symetric Noise $\rho_{-1} = \rho_{+1} = 0.4$ | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| **Datasets** | **Co-teaching+** | **Phuber** | **Bey.Ima.** | **PA** | **AdaStream** | **OLNLS** |
| breast | $.442 \pm .116$ | $.765 \pm .073$ | $.558 \pm .150$ | $.605 \pm .017$ | $.518 \pm .136$ | $\mathbf{.892 \pm .026}$ |
| dermatology | $.460 \pm .185$ | $.615 \pm .151$ | $.466 \pm .190$ | $.541 \pm .016$ | $.499 \pm .171$ | $\mathbf{.721 \pm .010}$ |
| wdbc | $.473 \pm .126$ | $.473 \pm .126$ | $.424 \pm .102$ | $.536 \pm .013$ | $.530 \pm .080$ | $\mathbf{.671 \pm .023}$ |
| diabetes | $.469 \pm .144$ | $.470 \pm .149$ | $.424 \pm .079$ | $.529 \pm .012$ | $.480 \pm .066$ | $\mathbf{.598 \pm .022}$ |
| german | $.449 \pm .172$ | $.539 \pm .195$ | $.606 \pm .155$ | $.509 \pm .013$ | $.377 \pm .094$ | $\mathbf{.694 \pm .006}$ |
| contraceptive | $.456 \pm .059$ | $.426 \pm .000$ | $.493 \pm .080$ | $.600 \pm .001$ | $.507 \pm .021$ | $\mathbf{.787 \pm .023}$ |
| splice | $.609 \pm .187$ | $.656 \pm .207$ | $.478 \pm .219$ | $.597 \pm .004$ | $.548 \pm .071$ | $\mathbf{.778 \pm .001}$ |
| mushroom | $.565 \pm .167$ | $.816 \pm .024$ | $.465 \pm .133$ | $.580 \pm .003$ | $.533 \pm .079$ | $\mathbf{.833 \pm .020}$ |
| marketing | $.691 \pm .254$ | $.794 \pm .018$ | $.699 \pm .229$ | $.562 \pm .002$ | $.517 \pm .121$ | $\mathbf{.818 \pm .005}$ |
| hapt | $.658 \pm .258$ | $.550 \pm .195$ | $.567 \pm .326$ | $.599 \pm .001$ | $.649 \pm .058$ | $\mathbf{.773 \pm .008}$ |
| ring | $.478 \pm .062$ | $.567 \pm .089$ | $.544 \pm .099$ | $.537 \pm .005$ | $.490 \pm .052$ | $\mathbf{.656 \pm .005}$ |
| magic04 | $.331 \pm .163$ | $.529 \pm .146$ | $.426 \pm .203$ | $.600 \pm .002$ | $.642 \pm .008$ | $\mathbf{.859 \pm .015}$ |
| a8a | $.657 \pm .098$ | $.689 \pm .009$ | $.467 \pm .232$ | $.574 \pm .002$ | $.687 \pm .004$ | $\mathbf{.696 \pm .019}$ |
| *loss/win* | 0/13 | 0/13 | 0/13 | 0/13 | 0/13 | $\mathbf{0/65}^*$ |
| *p*-value | .0001 | .0001 | .0001 | .0001 | .0001 | $--$ |
| *F*-rank | 2.3462 | 3.8077 | 2.3077 | 3.6154 | 2.9231 | 6.0000 |

Table 5: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR $\pm$ Standard Variance) for 13 datasets in the case of Symmetric Noisy labeled data. $\bullet$ indicates the cases that our method loses the comparison. * showing the total number of wins and losses for OLNLS. OOT: Out Of Time (24 hours).
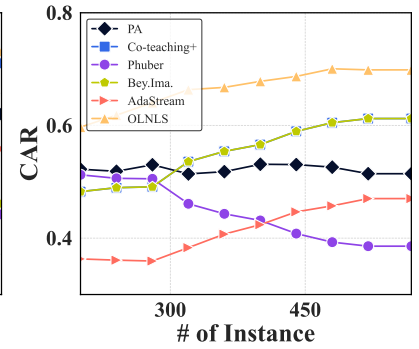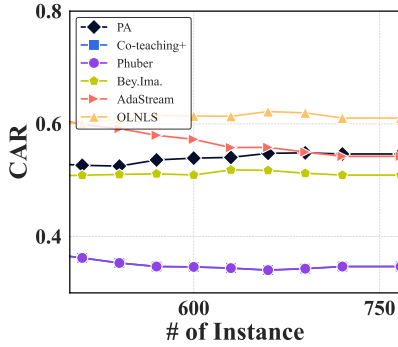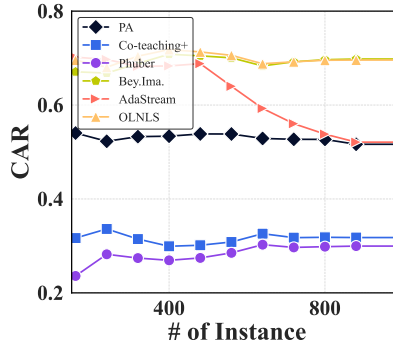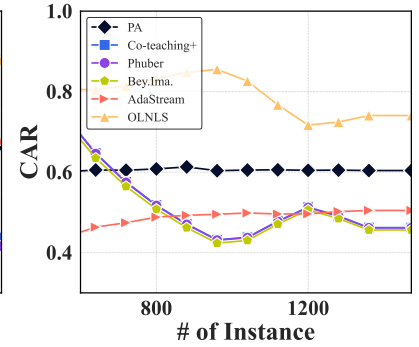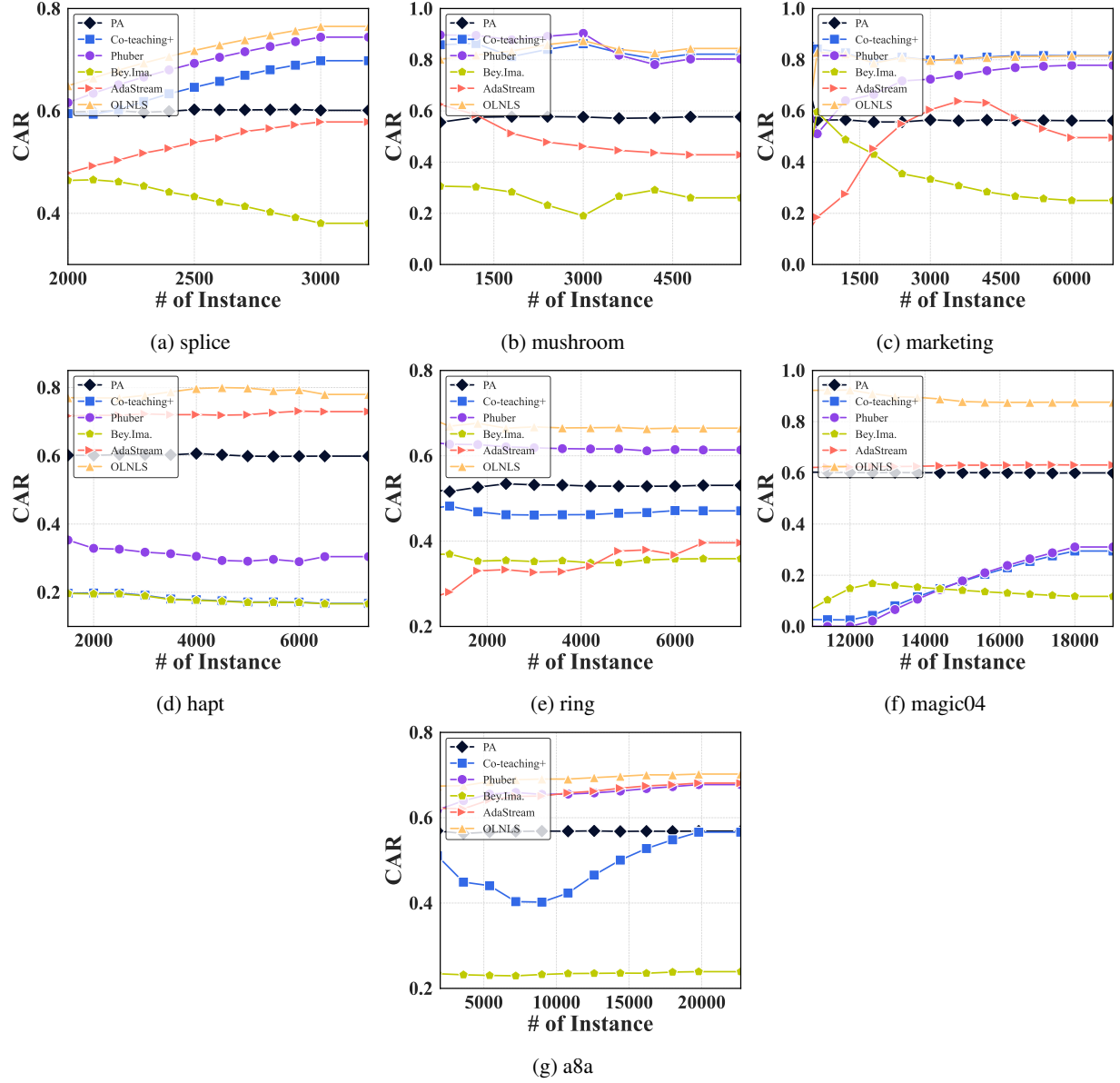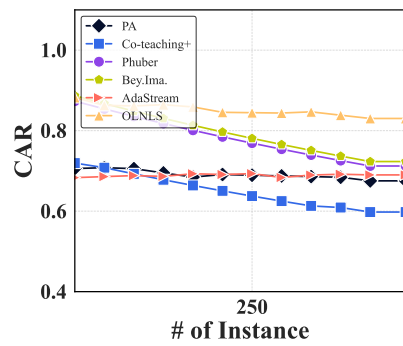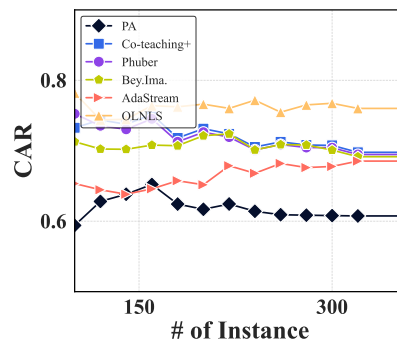
| | ASymetric Noise $\rho_{-1} = 0.3, \rho_{+1} = 0.5$ | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| **Datasets** | **Co-teaching+** | **Phuber** | **Bey.Ima.** | **PA** | **AdaStream** | **OLNLS** |
| breast | $.479 \pm .153$ | $.705 \pm .004$ | $.512 \pm .170$ | $.658 \pm .010$ | $.657 \pm .098$ | $\mathbf{.812 \pm .020}$ |
| dermatology | $.598 \pm .189$ | $.615 \pm .151$ | $.539 \pm .186$ | $.607 \pm .009$ | $.614 \pm .152$ | $\mathbf{.749 \pm .018}$ |
| wdbc | $.473 \pm .126$ | $.371 \pm .000$ | $.536 \pm .134$ | $.560 \pm .005$ | $.568 \pm .073$ | $\mathbf{.677 \pm .027}$ |
| diabetes | $.403 \pm .125$ | $.470 \pm .149$ | $.578 \pm .116\bullet$ | $.502 \pm .017$ | $.349 \pm .000$ | $.510 \pm .021$ |
| german | $.572 \pm .164\bullet$ | $.539 \pm .195$ | $.503 \pm .170$ | $.486 \pm .017$ | $.417 \pm .096$ | $.556 \pm .074$ |
| contraceptive | $.514 \pm .072$ | $.485 \pm .072$ | $.507 \pm .067$ | $.614 \pm .001$ | $.513 \pm .012$ | $\mathbf{.730 \pm .044}$ |
| splice | $.628 \pm .154$ | $.759 \pm .000$ | $.564 \pm .200$ | $.649 \pm .004$ | $.480 \pm .092$ | $\mathbf{.759 \pm .000}$ |
| mushroom | $.529 \pm .096$ | $.670 \pm .059$ | $.505 \pm .103$ | $.604 \pm .005$ | $.513 \pm .068$ | $\mathbf{.771 \pm .008}$ |
| marketing | $.507 \pm .277$ | $.816 \pm .002$ | $.553 \pm .298$ | $.629 \pm .003$ | $.488 \pm .085$ | $\mathbf{.826 \pm .002}$ |
| hapt | $.700 \pm .266$ | $.746 \pm .028$ | $.439 \pm .322$ | $.666 \pm .001$ | $.762 \pm .028$ | $\mathbf{.818 \pm .007}$ |
| ring | $.501 \pm .051$ | $.522 \pm .109$ | $.553 \pm .114\bullet$ | $.503 \pm .006$ | $.462 \pm .021$ | $.547 \pm .018$ |
| magic04 | $.545 \pm .145$ | $.679 \pm .045$ | $.415 \pm .116$ | $.631 \pm .002$ | $.642 \pm .007$ | $\mathbf{.811 \pm .022}$ |
| a8a | $.727 \pm .028$ | $.706 \pm .009$ | $.460 \pm .238$ | $.641 \pm .003$ | $.716 \pm .005$ | $\mathbf{.757 \pm .000}$ |
| *loss/win* | 1/12 | 0/13 | 2/11 | 0/13 | 0/13 | $\mathbf{3/62}^*$ |
| *p*-value | .0002 | .0013 | .0009 | .0001 | .0001 | $--$ |
| *F*-rank | 2.8462 | 3.8846 | 2.4615 | 3.3846 | 2.6923 | 5.7308 |

Table 6: The comparison results on cumulative accuracy rates. We repeated the experiment 5 times for each dataset, averaged the cumulative accuracy rate (CAR), and calculated the standard variance of the 5 times values. Experimental results (CAR $\pm$ Standard Variance) for 13 datasets in the case of Symmetric Noisy labeled data. ● indicates the cases that our method loses the comparison. * showing the total number of wins and losses for OLNLS. OOT: Out Of Time (24 hours)

| | Flip Noise rate = $0.4$ | | | | | |
|---|---|---|---|---|---|---|
| **Datasets** | **Co-teaching+** | **Phuber** | **Bey.Ima.** | **PA** | **AdaStream** | **OLNLS** |
| breast | $.492 \pm .094$ | $.692 \pm .070$ | $.551 \pm .194$ | $.609 \pm .009$ | $.508 \pm .127$ | $\mathbf{.896 \pm .052}$ |
| dermatology | $.530 \pm .182$ | $.464 \pm .185$ | $.574 \pm .153$ | $.563 \pm .028$ | $.538 \pm .186$ | $\mathbf{.745 \pm .019}$ |
| wdbc | $.576 \pm .102$ | $.525 \pm .126$ | $.475 \pm .126$ | $.533 \pm .015$ | $.576 \pm .063$ | $\mathbf{.673 \pm .023}$ |
| diabetes | $.410 \pm .123$ | $.470 \pm .149$ | $.500 \pm .130$ | $.540 \pm .013$ | $.406 \pm .071$ | $\mathbf{.579 \pm .018}$ |
| german | $.484 \pm .179$ | $.380 \pm .160$ | $.557 \pm .170$ | $.527 \pm .014$ | $.335 \pm .070$ | $\mathbf{.677 \pm .024}$ |
| contraceptive | $.539 \pm .103$ | $.456 \pm .059$ | $.485 \pm .072$ | $.600 \pm .001$ | $.500 \pm .023$ | $\mathbf{.790 \pm .028}$ |
| splice | $.483 \pm .202$ | $.345 \pm .207$ | $.626 \pm .196$ | $.596 \pm .005$ | $.488 \pm .086$ | $\mathbf{.779 \pm .001}$ |
| mushroom | $.437 \pm .092$ | $.667 \pm .178$ | $.526 \pm .112$ | $.582 \pm .003$ | $.499 \pm .060$ | $\mathbf{.827 \pm .006}$ |
| marketing | $.502 \pm .276$ | $.804 \pm .012$ | $.673 \pm .247$ | $.564 \pm .004$ | $.447 \pm .105$ | $\mathbf{.815 \pm .006}$ |
| hapt | $.688 \pm .260$ | $.766 \pm .088$ | $.757 \pm .124$ | $.597 \pm .001$ | $.713 \pm .044$ | $\mathbf{.772 \pm .014}$ |
| ring | $.484 \pm .029$ | $.388 \pm .000$ | $.512 \pm .054$ | $.535 \pm .005$ | $.508 \pm .042$ | $\mathbf{.660 \pm .002}$ |
| magic04 | $.652 \pm .058$ | $.689 \pm .053$ | $.546 \pm .124$ | $.599 \pm .002$ | $.642 \pm .008$ | $\mathbf{.869 \pm .014}$ |
| a8a | $.636 \pm .198$ | $.678 \pm .012$ | $.460 \pm .191$ | $.576 \pm .003$ | $.684 \pm .004$ | $\mathbf{.711 \pm .002}$ |
| *loss/win* | $0/13$ | $0/13$ | $0/13$ | $0/13$ | $0/13$ | $\mathbf{0/65}^*$ |
| *p*-value | $.0001$ | $.0001$ | $.0001$ | $.0001$ | $.0001$ | $--$ |
| *F*-rank | $2.5000$ | $3.0769$ | $3.2308$ | $3.5385$ | $2.6538$ | $6.000$ |



(a) breast    (b) dermatology    (c) wdbc

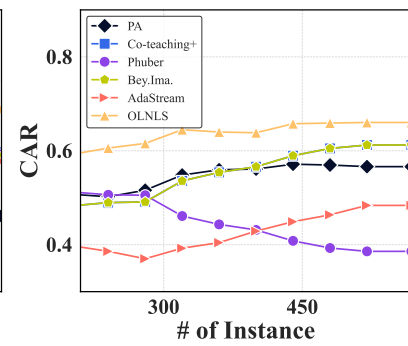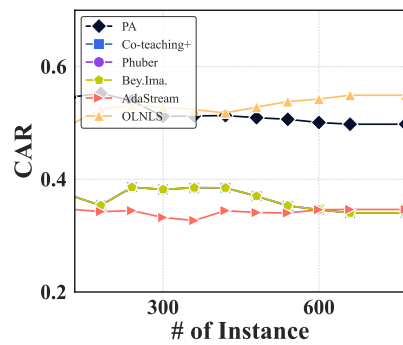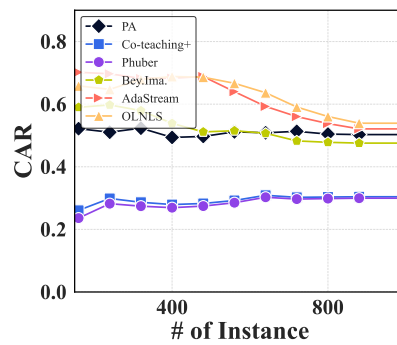(d) diabetes    (e) german    (f) contraceptive

Figure 4: The cumulative accuracy rate (CAR) trends of Co-teaching+, Phuber, BeyondImages, PA, AdaStream, and OLNLS in all 13 symmetric noisy labeled data streams.
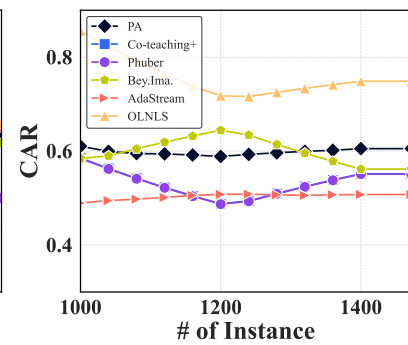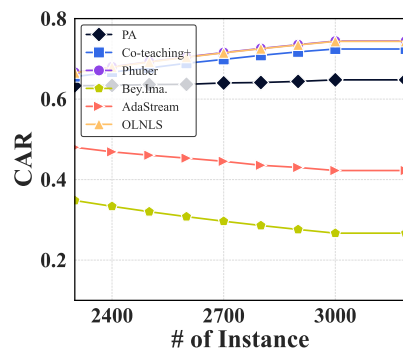
(a) breast      (b) dermatology      (c) wdbc

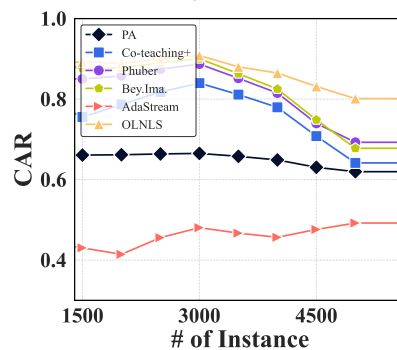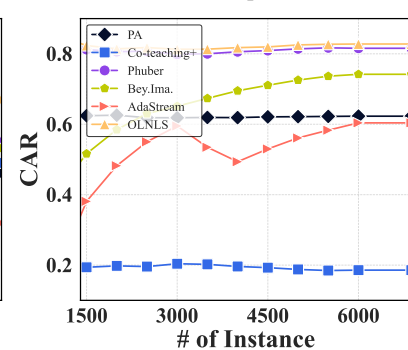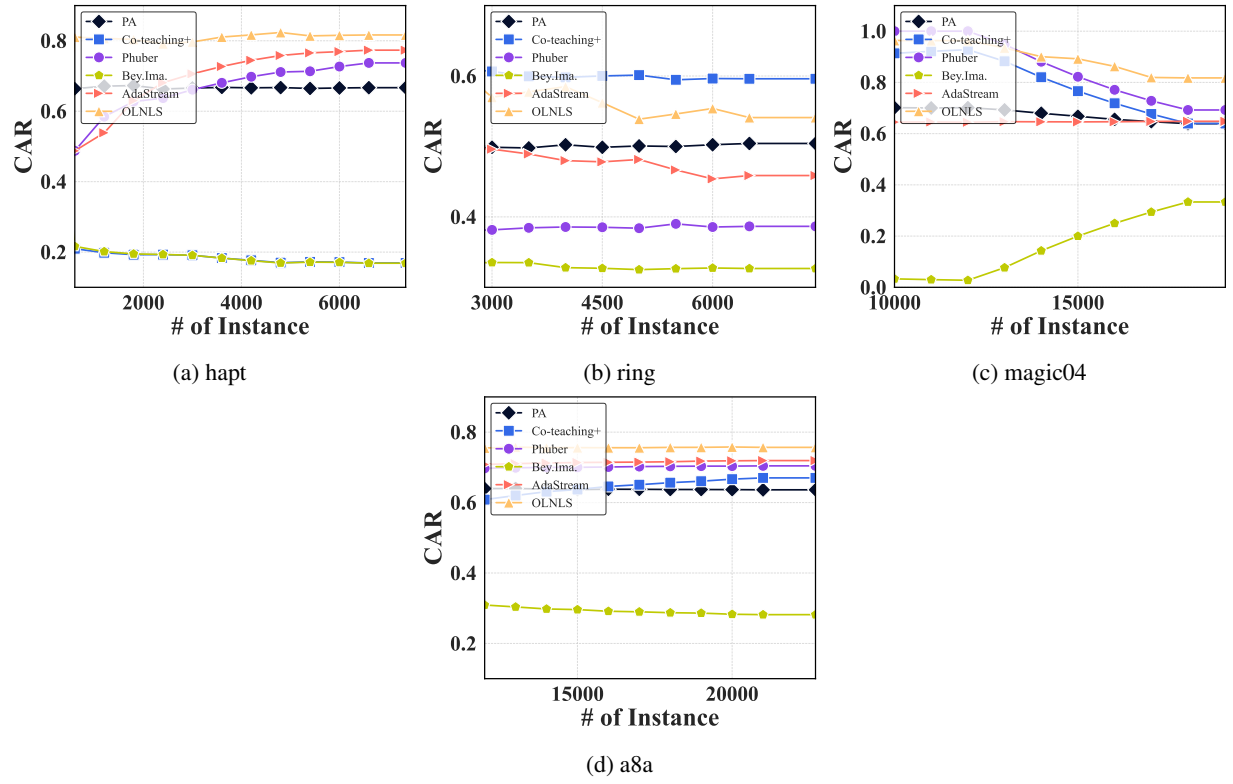(d) diabetes      (e) german      (f) contraceptive
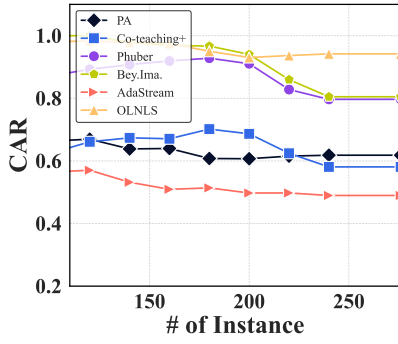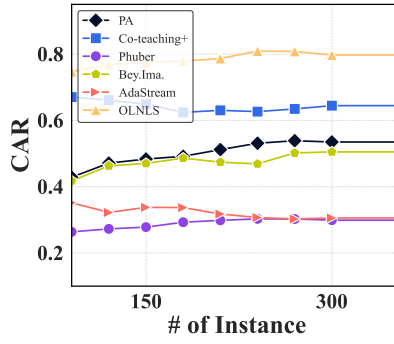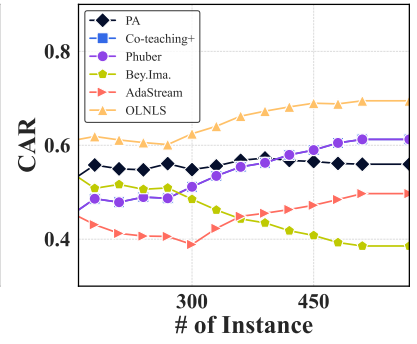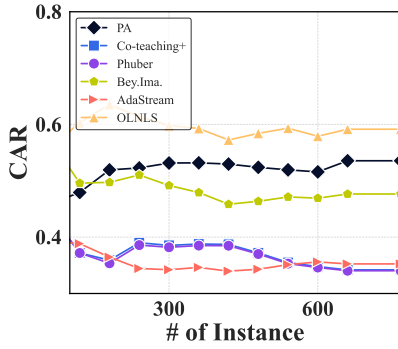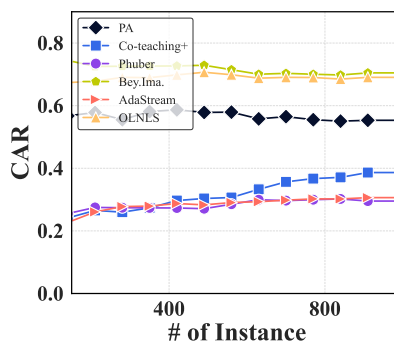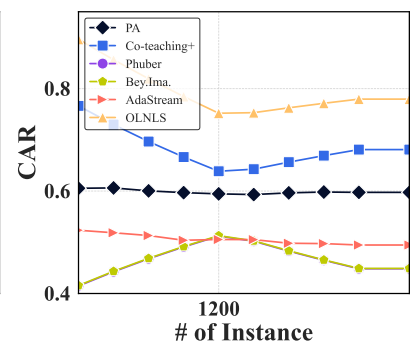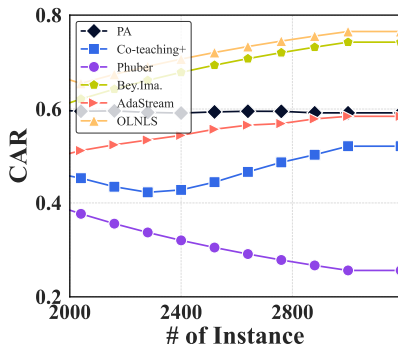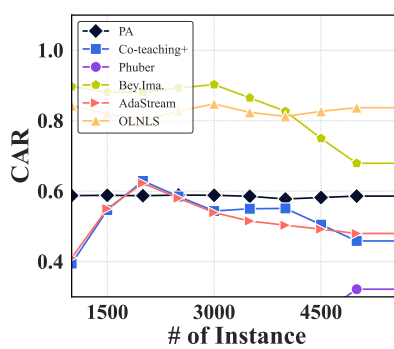
(g) splice      (h) mushroom      (i) marketing

(a) hapt

(b) ring

(c) magic04

(d) a8a

Figure 5: The cumulative accuracy rate (CAR) trends of Co-teaching+, Phuber, BeyondImages, PA, AdaStream, and OLNLS in all 13 asymmetric noisy labeled data streams.
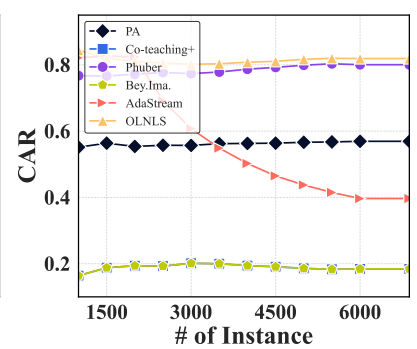
(a) breast      (b) dermatology      (c) wdbc

(d) diabetes      (e) german      (f) contraceptive

(g) splice      (h) mushroom      (i) marketing
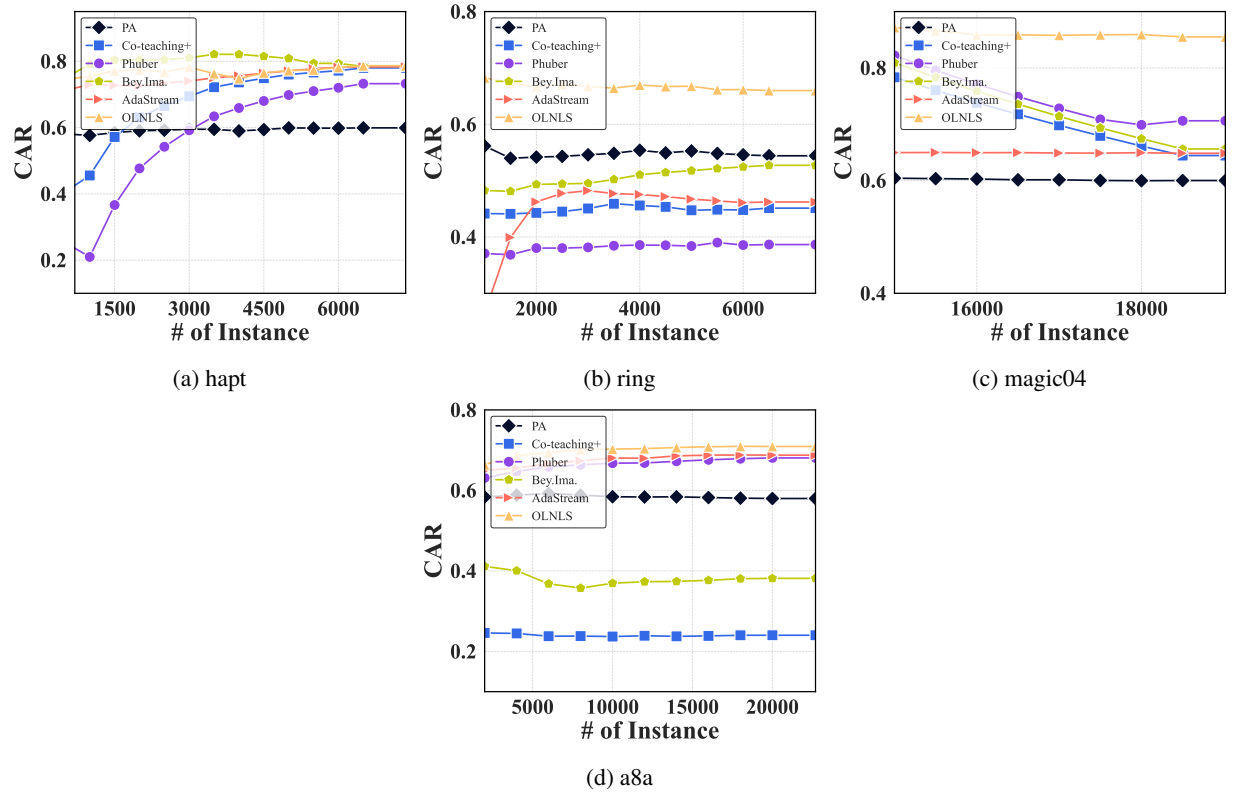
(a) hapt

(b) ring

(c) magic04

(d) a8a

Figure 6: The cumulative accuracy rate (CAR) trends of Co-teaching+, Phuber, BeyondImages, PA, AdaStream, and OLNLS in all 13 flipped noisy labeled data streams.