

556 A Further Exposition on Usage of AGTRs

557 We will use this first section of our appendix to provide further exposition as to how we envision
558 AGTRs to be used, as well as their constraints. To be clear, our proposed AGTR does not alleviate the
559 need for labeled data, or avoid all possible biases. Indeed by constructing a mechanism that produces
560 AGTRs, assuming it does not produce a complete graph, there must necessarily be some form of
561 bias that exists via the selection of the AGTR production itself. Not only is labeled data needed to
562 produce the initial model being evaluated, but labeled data is necessary to produce the estimate of the
563 error rate $\hat{\epsilon}$ of the AGTR construction. These are important nuances to keep in mind through usage.

564 Our belief is that there are many domains where constructing an AGTR, with confidently small $\hat{\epsilon}$, is
565 possible. Our example application has been for malware analysis as it is the area of our expertise.
566 For example we see potential application areas in chemistry, physics, optics, and others where
567 obtained experimental ground truth results is highly expensive and time consuming, but sophisticated
568 simulations may be able to group different scenarios together with high fidelity. That is that the
569 simulation's result may not produce the exactly correct prediction, but could be used to group different
570 scenarios being evaluated into groups highly probably to have outcomes of the same fundamental
571 nature. This is extrapolation on our part since we do not have the requisite background to delve
572 deeply into these other areas, but captures our primary hypothesis: many domains have significant
573 domain knowledge that could produce viable AGTRs.

574 Once an AGTR is produced, it provides a valuable means of improving bench-marking, confidence in
575 results, and testing of nuanced changes. Designing good benchmarks if of critical importance, and
576 an AGTR can be an effective litmus test for any benchmark creation. It allows one to cheaply use a
577 much larger set of unlabeled data to determine if the benchmark itself has become the subject of over-
578 fitting. We provided an example of this with the AVClass tool in our paper, where our bounds indicate
579 AVClass's day may be over-fit because their empirical recall is higher than our upper bounded recall.
580 This does not necessarily mean the results are actually over-fit, because the datasets evaluated are
581 not the same (in AVClass' cases, not available to us). But it does provide a strong indicator that
582 something is amiss, and provides the created of a benchmark the signal needed to more thoroughly
583 and rigorously investigate their data.

584 This utility of AGTRs in benchmarking is simultaneously macro and micro in nature. The prior
585 paragraph has described the macro nature, that we can perform litmus testing against a created
586 benchmark to detect overfitting against a larger unlabeled corpus. The next level down that one would
587 natural want to pursue is to compare two different models using an AGTR, but we must caution
588 the reader into doing such a comparison. When two models are being evaluated and of different
589 mechanical natures, they will naturally correlate to different degrees with the biases of the AGTR
590 construction itself. This means one method could appear to have better predicted bounds by our AGTR
591 by only being successful on the data for which the AGTR itself is able to operate, while potentially
592 being errant on all other data and would not be realized. Any other more reasonable model, that
593 simply lacks the bias to the AGTR in use, would then appear worse when it may in fact be superior.

594 This is the reason why we recommend to only use AGTR to compare models that are fundamentally
595 very similar, like changing small parameters of a single approach. This is reasonable because the
596 underlying mechanism of prediction remains the same, and allows us to better quantify the impacts
597 of minute changes that would be difficult to estimate on small static corpora. This assumes that the
598 parameters themselves do not have an out-sized impact on the fundamental nature, which we believe
599 is reasonable but should be made explicit.

600 This highlights another limitation of our approach, which is the number of epistemic uncertainties
601 involved in its use. It will depend on the beliefs and confidence of the domain experts that the models
602 being compared are "similar enough" that using an AGTR to compare changes is valid, and that the
603 error rate bound $\hat{\epsilon}$ is sufficient. While more concrete answers to these are desirable, we recognize
604 them as important items of future work. As we have developed the AGTR approach thus far, it
605 provides a powerful means of better leveraging unlabelled data in the evaluation of a ML model,
606 which we argue is additive to the tools in use today.

607 B Applications of Approximate Ground Truth Refinements

608 As long as an AGTR can be constructed from a dataset, that dataset can be used for partial evaluation
609 of a clustering algorithm or multiclass classifier - even if the dataset does not have reference labels.
610 This is notable because it allows larger, more representative datasets without reference labels to be
611 used during the evaluation process. The ability to compute bounds on precision, recall, and accuracy
612 without reference labels is valuable in its own right, but we have found two additional ways in which
613 AGTRs can be used to great effect. First, the bounds from an AGTR can be used as a litmus test for
614 detecting biased evaluation results produced using a substandard reference dataset. Second, AGTRs
615 can be used for evaluating modifications to a clustering algorithm or multiclass classifier.

616 B.1 Testing Suspect Evaluation Results

617 We have found that AGTRs can be used to detect misleading results produced from low quality
618 reference datasets. The following five steps describe how to test suspect evaluation results:

- 619 1. Compute the precision, recall, and/or accuracy of a clustering algorithm or a multiclass
620 classifier using a substandard reference dataset.
- 621 2. Obtain C by applying the clustering algorithm or multiclass classifier to a large, diverse
622 dataset.
- 623 3. Construct an AGTR \hat{R} from the dataset in step 2.
- 624 4. Select a value of $\hat{\epsilon}$ that is believed to be greater than the number of errors in \hat{R} .
- 625 5. Compute $Precision(C, \hat{R}) - \frac{\hat{\epsilon}}{m}$ and $Recall(C, \hat{R}) + \frac{\hat{\epsilon}}{m}$. Test whether all bounds hold for the
626 evaluation results found during step 1.

627 Note that if testing a multiclass classifier, C is obtained during step 2 by using the classifier to predict
628 the class label for each data point in the dataset and then clustering all data points that share the same
629 predicted class label. Subsection 2.4 describes the qualities an ideal AGTR should possess for best
630 evaluation results. Refer to subsection 2.5 for discussion about how to select an appropriate value of
631 $\hat{\epsilon}$ during step 4.

632 Because the dataset used to construct the AGTR during step 3 is larger and more diverse than the
633 reference dataset from step 1, we assume that the AGTR dataset is a better exemplar of the overall
634 problem space. Although the metric bounds found using the AGTR dataset do not necessarily hold
635 for the reference dataset, it would still be irregular for any evaluation results from the reference
636 dataset to violate them. Therefore, if $Precision(C, \hat{R}) - \frac{\hat{\epsilon}}{m}$ is greater than the precision computed
637 using the reference dataset, or if $Recall(C, \hat{R}) + \frac{\hat{\epsilon}}{m}$ is less than the recall or accuracy computed
638 using the reference dataset, it is reasonable to conclude that the evaluation results found using the reference
639 dataset are misleading.

640 B.2 Comparing Similar Models

641 An important component of cluster and classifier evaluation is the ability to compare models against
642 each other. Suppose we wish to determine which of two clustering algorithms has a higher precision,
643 but we do not have access to a satisfactory reference dataset. We use the two models to predict
644 clusterings C_1 and C_2 from an unlabeled dataset and we construct an AGTR \hat{R} from the same dataset.
645 The precision lower bounds of the two clustering algorithms are given by $Precision(C_1, \hat{R}) - \frac{\hat{\epsilon}}{m}$
646 and $Precision(C_2, \hat{R}) - \frac{\hat{\epsilon}}{m}$. Although we can apply Theorem 5 to show $Precision(C_1, \hat{R}) - \frac{\hat{\epsilon}}{m} \leq$
647 $Precision(C_1, D)$ and $Precision(C_2, \hat{R}) - \frac{\hat{\epsilon}}{m} \leq Precision(C_2, D)$, we cannot prove any relationship
648 between $Precision(C_1, D)$ and $Precision(C_2, D)$.

649 Unfortunately, evaluation metric bounds cannot be used to provably determine whether one clustering
650 algorithm or multiclass classifier has a higher precision, recall, or accuracy than another. However, in
651 specific cases higher evaluation metric bounds may indicate that one model has a higher performance
652 than another. Two conditions must be met in order for this approach to be used. First, the clustering
653 algorithms or classifiers being compared must be intrinsically similar, such as two different versions
654 of the same classifier. Second, one of the clustering algorithms or classifiers must be tested to ensure

655 that changes in performance are strongly correlated to changes in evaluation metric bounds. The
656 steps of the test are as follows:

- 657 1. Obtain C by applying the clustering algorithm or multiclass classifier to a large, diverse
658 dataset.
- 659 2. Construct an AGTR \hat{R} from the dataset in step 1.
- 660 3. Incrementally shuffle the cluster membership of each data point in C .
- 661 4. Compute $Precision(C, \hat{R}) - \frac{\epsilon}{m}$ and $Recall(C, \hat{R}) + \frac{\epsilon}{m}$ at regular intervals of the shuffle.
- 662 5. Compute correlation between shuffle percentage and the evaluation metric bounds. Test that
663 a strong negative correlation between the two exists.

664 Step 3 is performed randomly sampling a data point in C with no replacement, randomly selecting a
665 cluster in C weighted by the original distribution of cluster sizes, and then assigning the data point
666 to that cluster. The process repeats m times, where m is the number of points in the dataset, after
667 which the entire clustering has been shuffled. Because the dataset has high diversity, the probability
668 that a datapoint is randomly assigned to an incorrect cluster is far greater the probability that it is
669 randomly assigned the correct one. Therefore, shuffling C with this strategy is very likely to produce
670 predicted clusterings that are sequentially worse. Step 5 is performed by computing the Pearson
671 correlation between the shuffle percentage and the evaluation metric bounds. If a very strong negative
672 correlation exists between the two, we conclude that small modifications to a clustering algorithm or
673 classifier that improve its performance will be reflected by a higher evaluation metric bounds, and
674 modifications that lower its performance will be reflected by lower bounds.

675 C Malware Label and Dataset Challenges

676 The task of labeling malware data is expensive and error prone, more so than most standard ML
677 applications areas. This forces datasets to be either small and thus non-representative of the large
678 and diverse malware ecosystem, or large but noisily labeled or biased which again limits conclusions.
679 Because of these issues, we believe that malware family classification is an ideal problem space for
680 demonstrating the AGTR evaluation framework. We will quickly review the challenges in labeling
681 malware and how it has caused a lack of quality datasets.

682 Ideally a malware reference dataset would be constructed using *manual labeling* to determine the
683 malware family of each file. Manual analysis is not perfectly accurate, but the error rate is considered
684 negligible enough that labels obtained via manual analysis are considered to have ground truth
685 confidence [13]. A professional analyst can take a 10 hours or more to fully analyze a single
686 file [14, 15]. This level of analysis is not always needed to determine the family of a malware sample,
687 but it exemplifies the high human cost of manual labeling that prevents its use in all but very small
688 datasets [1]. One approach to mitigate this is *cluster labeling*, where the dataset is clustered and an
689 exemplar from each cluster is manually labeled. This strategy is highly reliant on the precision of the
690 algorithm used to cluster the reference dataset, which is often custom-made [16, 17]. Furthermore,
691 the scalability of cluster labeling is still limited due to the requirement of manual analysis.

692 For these reasons larger malware reference datasets tend to use *antivirus labeling*, where an antivirus
693 engine is used to label a corpus. This is easy to implement at scale, but has significant quality
694 issues [18]. Antivirus labels are frequently incomplete, inconsistent or incorrect [19, 20]. Antivirus
695 signatures do not always contain family information [13], and different antivirus engines disagree on
696 the names of malware families [9]. Labels from an antivirus engine can take almost a year to stabilize
697 [18], creating a necessary lag time between data occurrence and label inference. These issues can be
698 partially mitigated by using *antivirus majority voting*, but aggregating the results of multiple antivirus
699 engines produces unlabeled files due to lack of consensus. This then biases the final dataset to only
700 the “easy” samples that are the least interesting and already known by current tools [2].

701 The aforementioned labeling issues have a significant impact on the datasets that are available. The
702 largest datasets used in malware classifier and clustering research range from a hundred thousand
703 [13, 16] up to one million samples[21], but are private corpora held by corporations that can afford
704 the construction cost and do not want to give away a competitive advantage¹. Since the data is private

¹There are also legal concerns for sharing benign applications, but our discussion is focused solely on malware.

705 the validation of the labeling can not be replicated or investigated, and in most cases the number
 706 of families is not fully specified[13, 16]. The vast majority of publicly available datasets that have
 707 been used are less than 12,000 samples in size[17, 22–26]. Of these MalGenome is the only fully
 708 manually labeled corpus, but also the smallest with only 49 families and 1,260 files[27, 28]. The
 709 small size and low diversity of these corpora makes it difficult to make generalizable conclusions
 710 about the quality of a malware clustering algorithm or classifier. The largest reference dataset in
 711 general use is VX Heavens with 271,092 samples, but it contains very old malware, is labeled using
 712 only a single antivirus engine, and the history and composition of the dataset is poorly documented
 713 [29]. Additional statistics about notable malware reference datasets are listed in Appendix E.

714 Given the issues in curating these datasets, we have identified three common three common traits
 715 that negatively impact evaluation. All of the datasets we just referenced possess one or more of these
 716 undesirable traits.

717 1) *Reference labels without ground truth confidence*: Because producing ground truth reference labels
 718 for a large corpus of malware samples is infeasible, it is common practice to use malware reference
 719 labels without ground truth confidence [1]. Many prior malware classifiers have been evaluated using
 720 non-ground truth reference labels without the quality of those labels having been assessed [9], which
 721 may result in overoptimistic or misleading results [2].

722 2) *Insufficient size or diversity*: Due to the enormous number of families in existence, malware
 723 reference datasets must be both large and diverse in order to be representative of the malware
 724 ecosystem. Using a reference dataset with a small number, or imbalanced set, of families lowers the
 725 significance of evaluation results [2].

726 3) *Outdated malware samples*: The ecosystem of malware is constantly changing as categories of
 727 malware, malware families, and other tradecraft rise to prominence or fall out of favor. Evaluating a
 728 clustering algorithm or multiclass classifier using a dataset of outdated malware may produce results
 729 that do not translate to present day malware. All of the discussed datasets contain malware samples
 730 from 2015 and earlier, failing to represent the last half decade of malware development.

731 D Proof Details

732 **Theorem 1.** $Precision(C, R) \leq Precision(C, D)$

Proof of Theorem 1. Suppose some i s.t. $1 \leq i \leq c$. Since R is a refinement of D , $\exists D_j \in D$ s.t.
 $R_{f'(i)} \subseteq D_j$. Because $R_{f'(i)} \subseteq D_j$, it must be that $|C_i \cap R_{f'(i)}| \leq |C_i \cap D_j|$. By definition, $f(i) =$
 $\operatorname{argmax}_j |C_i \cap D_j|$. Therefore, $|C_i \cap R_{f'(i)}| \leq |C_i \cap D_j| \leq |C_i \cap D_{f(i)}|$. We can simply sum over
 this inequality to obtain:

$$\frac{1}{m} \sum_{i=1}^c |C_i \cap R_{f'(i)}| \leq \frac{1}{m} \sum_{i=1}^c |C_i \cap D_{f(i)}|$$

733 By Definition 1, $Precision(C, R) \leq Precision(C, D)$. □

734 **Theorem 2.** $Recall(C, R) \geq Recall(C, D)$

Proof of Theorem 2. Suppose some j s.t. $1 \leq j \leq d$. Because R is a refinement of D , $\exists Q_j =$
 $\{Q_{j\ell}\}_{1 \leq \ell \leq q_j}$ s.t. $Q_{j\ell} \in R$ and $D_j = \bigcup_{\ell=1}^{q_j} Q_{j\ell}$. We can say that $|D_j| = \sum_{\ell=1}^{q_j} |Q_{j\ell}|$ and furthermore
 that $|C_{g(j)} \cap D_j| = \sum_{\ell=1}^{q_j} |C_{g(j)} \cap Q_{j\ell}|$. We know that $\forall Q_{j\ell} \in Q_j, \exists R_k \in R$ s.t. $Q_{j\ell} = R_k$. By
 definition $g'(k) = \operatorname{argmax}_i |C_i \cap R_k|$, and $R_k = Q_{j\ell}$, so $|C_{g'(k)} \cap R_k| \geq |C_{g(j)} \cap Q_{j\ell}|$. Since by
 Property 1 the sets R_k are in bijection with the sets $Q_{j\ell}$, we can sum over this inequality to obtain:

$$\frac{1}{m} \sum_{k=1}^r |C_{g'(k)} \cap R_k| \geq \frac{1}{m} \sum_{j=1}^d \sum_{\ell=1}^{q_j} |C_{g(j)} \cap Q_{j\ell}|$$

$$\frac{1}{m} \sum_{j=1}^d \sum_{\ell=1}^{q_j} |C_{g(j)} \cap Q_{j\ell}| = \frac{1}{m} \sum_{j=1}^d |C_{g(j)} \cap D_j|$$

735 By Definition 2, $\text{Recall}(C, R) \geq \text{Recall}(C, D)$. \square

736 **Corollary 2.1.** $\text{Recall}(C, R) \geq \text{Accuracy}(C, D)$

737 *Proof of Corollary 2.1.* Using Theorem 2, $\text{Recall}(C, R) \geq \text{Recall}(C, D)$. It must be the case that
738 $\text{Recall}(C, D)$ where $g(j) = \arg\max_i |C_i \cap D_j| \geq \text{Recall}(C, D)$ where g is the identity function. Using
739 Definition 3, $\text{Recall}(C, D) = \text{Accuracy}(C, D)$ when g is the identity function. Therefore, $\text{Recall}(C, R)$
740 $\geq \text{Accuracy}(C, D)$. \square

741 **Theorem 3.** $|\text{Precision}(C, S) - \text{Precision}(C, \hat{S})| \leq \frac{1}{m}$

742 *Proof of Theorem 3.* Let $S = \{S_t\}_{1 \leq t \leq s}$ be an arbitrary partition of M . Let the label translation
743 function $f : \{1 \dots c\} \mapsto \{1 \dots s\}$ be defined as $f(i) = \arg\max_t |C_i \cap S_t|$. Suppose some $M_n \in M$,
744 some $S_x \in S$ s.t. $M_n \in S_x$, and some S_y s.t. $S_y \in S$ or $S_y = \{\emptyset\}$. Furthermore, suppose
745 some $C_a, C_b \in C$ s.t. $M_n \in C_a$ and $b = f(y)$. Let $\hat{S} = \{\hat{S}_t\}_{1 \leq t \leq s}$ be a clustering identical to
746 S except for one cluster label change, which is given by $\hat{S}_x = S_x - \{M_n\}$ and $\hat{S}_y = S_y \cup \{M_n\}$.
747 Let the function $\hat{f} : \{1 \dots c\} \mapsto \{1 \dots \hat{s}\}$ be defined as $\hat{f}(i) = \arg\max_{\hat{a}} |C_i \cap \hat{S}_{\hat{a}}|$. At minimum
748 $|C_a \cap \hat{S}_{\hat{f}(a)}| = |C_a \cap S_{f(a)}| - 1$ and at maximum $|\hat{S}_{\hat{f}(a)}| = |C_a \cap S_{f(a)}|$. Similarly, at minimum
749 $|C_b \cap \hat{S}_{\hat{f}(b)}| = |C_b \cap S_{f(b)}|$ and at maximum $|C_b \cap \hat{S}_{\hat{f}(b)}| = |C_b \cap S_{f(b)}| + 1$. Finally, because
750 all elements in M are unique and C partitions M , $\forall C_i \in C$ s.t. $i \neq a$ and $i \neq b$, $|C_i \cap \hat{S}_{\hat{f}(i)}| =$
751 $|C_i \cap S_{f(i)}|$. Therefore, $\sum_{t=1}^s |C_{g(t)} \cap S_t| - 1 \leq \sum_{t=1}^{\hat{s}} |C_{\hat{g}(t)} \cap \hat{S}_t| \leq \sum_{t=1}^s |C_{g(t)} \cap S_t| + 1$. By Definition
752 1, $\text{Precision}(C, S) - \frac{1}{m} \leq \text{Precision}(C, \hat{S}) \leq \text{Precision}(C, S) + \frac{1}{m}$. We write this as $|\text{Precision}(C, S)$
753 $- \text{Precision}(C, \hat{S})| \leq \frac{1}{m}$. \square

754 **Corollary 3.1.** $|\text{Precision}(C, R) - \text{Precision}(C, \hat{R})| \leq \frac{\epsilon}{m}$

755 *Proof of Corollary 3.1.* By Theorem 3, $|\text{Precision}(C, S) - \text{Precision}(C, \hat{S})| \leq \frac{1}{m}$ for some arbitrary
756 clustering S and a second clustering \hat{S} that is equivalent to S but with a single data point belonging
757 to a different cluster. Given a GTR R and a corresponding AGTR \hat{R} , we can sequentially change the
758 cluster membership of ϵ data points in R to obtain \hat{R} . At each step the precision value can change by
759 at most $\pm \frac{1}{m}$. Therefore, $|\text{Precision}(C, R) - \text{Precision}(C, \hat{R})| \leq \frac{\epsilon}{m}$. \square

760 **Corollary 3.2.** $|\text{Precision}(C, D) - \text{Precision}(C, \hat{R})| \leq \frac{\delta}{m}$

761 *Proof of Corollary 3.2.* By Theorem 3, $|\text{Precision}(C, S) - \text{Precision}(C, \hat{S})| \leq \frac{1}{m}$ for some arbitrary
762 clustering S and a second clustering \hat{S} that is equivalent to S but with a single data point belonging
763 to a different cluster. Given a ground truth reference clustering D and a corresponding AGTR \hat{R} , we
764 can sequentially change the cluster membership of δ data points in D to obtain \hat{R} . At each step the
765 precision value can change by at most $\pm \frac{1}{m}$. Therefore, $|\text{Precision}(C, D) - \text{Precision}(C, \hat{R})| \leq \frac{\delta}{m}$. \square

766 **Theorem 4.** $|\text{Recall}(C, S) - \text{Recall}(C, \hat{S})| \leq \frac{1}{m}$

767 *Proof of Theorem 4.* Let $S = \{S_t\}_{1 \leq t \leq s}$ be an arbitrary partition of M . Let the label translation
768 function $g : \{1 \dots s\} \mapsto \{1 \dots c\}$ be defined as $g(t) = \operatorname{argmax}_i |C_i \cap S_t|$. Suppose some $M_n \in M$,
769 some $S_x \in S$ s.t. $M_n \in S_x$, and some S_y s.t. $S_y \in S$ or $S_y = \{\emptyset\}$. Let $\hat{S} = \{\hat{S}_t\}_{1 \leq t \leq \hat{s}}$ be a
770 clustering identical to S except for one cluster label change, which is given by $\hat{S}_x = S_x - \{M_n\}$ and
771 $\hat{S}_y = S_y \cup \{M_n\}$. Let the function $\hat{g} : \{1 \dots \hat{s}\} \mapsto \{1 \dots c\}$ be defined as $\hat{g}(t) = \operatorname{argmax}_i |C_i \cap \hat{S}_t|$. At
772 minimum $|C_{\hat{g}(x)} \cap \hat{S}_x| = |C_{g(x)} \cap S_x| - 1$ and at maximum $|C_{\hat{g}(x)} \cap \hat{S}_x| = |C_{g(x)} \cap S_x|$. Similarly,
773 at minimum $|C_{\hat{g}(y)} \cap \hat{S}_y| = |C_{g(y)} \cap S_y|$ and at maximum $|C_{\hat{g}(y)} \cap \hat{S}_y| = |C_{g(y)} \cap S_y| + 1$. Because
774 each other clusters in \hat{S} identical to some cluster in S , $\sum_{t=1}^{\hat{s}} |C_{\hat{g}(t)} \cap \hat{S}_t| - 1 \leq \sum_{t=1}^{\hat{s}} |C_{\hat{g}(t)} \cap \hat{S}_t| \leq$
775 $\sum_{t=1}^s |C_{g(t)} \cap S_t| + 1$. Using Definition 2, we obtain $\operatorname{Recall}(C, S) - \frac{1}{m} \leq \operatorname{Recall}(C, \hat{S}) \leq \operatorname{Recall}(C, S)$
776 $+ \frac{1}{m}$. We write this as $|\operatorname{Recall}(C, S) - \operatorname{Recall}(C, \hat{S})| \leq \frac{1}{m}$. \square

777 **Corollary 4.1.** $|\operatorname{Recall}(C, R) - \operatorname{Recall}(C, \hat{R})| \leq \frac{\epsilon}{m}$

778 *Proof of Corollary 4.1.* By Theorem 4, $|\operatorname{Recall}(C, S) - \operatorname{Recall}(C, \hat{S})| \leq \frac{1}{m}$ for some arbitrary
779 clustering S and a second clustering \hat{S} that is equivalent to S but with a single data point belonging
780 to a different cluster. Given a GTR R and a corresponding AGTR \hat{R} , we can sequentially change the
781 cluster membership of ϵ data points in R to obtain \hat{R} . At each step the recall value can change by at
782 most $\pm \frac{1}{m}$. Therefore, $|\operatorname{Recall}(C, R) - \operatorname{Recall}(C, \hat{R})| \leq \frac{\epsilon}{m}$. \square

783 **Corollary 4.2.** $|\operatorname{Recall}(C, D) - \operatorname{Recall}(C, \hat{R})| \leq \frac{\delta}{m}$

784 *Proof of Corollary 4.2.* By Theorem 4, $\operatorname{Recall}(C, S) - \frac{1}{m} \leq \operatorname{Recall}(C, \hat{S}) \leq \operatorname{Recall}(C, S) + \frac{1}{m}$
785 for some arbitrary clustering S and a second clustering \hat{S} that is equivalent to S but with a single
786 data point belonging to a different cluster. Given a ground truth reference label clustering D and a
787 corresponding AGTR \hat{R} , we can sequentially change the cluster membership of ϵ data points in D
788 to obtain \hat{R} . At each step the recall value can change by at most $\pm \frac{1}{m}$. Therefore, $|\operatorname{Recall}(C, D) -$
789 $\operatorname{Recall}(C, \hat{R})| \leq \frac{\delta}{m}$. \square

790 **Theorem 5.** If $\hat{e} \geq e$ then $\operatorname{Precision}(C, \hat{R}) - \frac{\epsilon}{m} \leq \operatorname{Precision}(C, D)$

791 *Proof of Theorem 5.* By Corollary 3.1, $|\operatorname{Precision}(C, R) - \operatorname{Precision}(C, \hat{R})| \leq \frac{\epsilon}{m}$. This can be
792 written as $\operatorname{Precision}(C, \hat{R}) - \frac{\epsilon}{m} \leq \operatorname{Precision}(C, R)$. By applying Theorem 1, $\operatorname{Precision}(C, \hat{R}) - \frac{\epsilon}{m} \leq$
793 $\operatorname{Precision}(C, \hat{R}) - \frac{\epsilon}{m} \leq \operatorname{Precision}(C, R) \leq \operatorname{Precision}(C, D)$. \square

794 **Theorem 6.** If $\hat{e} \geq e$ then $\operatorname{Recall}(C, \hat{R}) + \frac{\epsilon}{m} \geq \operatorname{Recall}(C, D)$

795 *Proof of Theorem 6.* By Corollary 4.1, $|\operatorname{Recall}(C, R) - \operatorname{Recall}(C, \hat{R})| \leq \frac{\epsilon}{m}$. This can be written as
796 $\operatorname{Recall}(C, \hat{R}) + \frac{\epsilon}{m} \geq \operatorname{Recall}(C, R)$. Therefore, by Theorem 2, $\operatorname{Recall}(C, \hat{R}) + \frac{\epsilon}{m} \geq \operatorname{Recall}(C, \hat{R}) +$
797 $\frac{\epsilon}{m} \geq \operatorname{Recall}(C, R) \geq \operatorname{Recall}(C, D)$. \square

798 **Corollary 6.1.** If $\hat{e} \geq e$ then $\operatorname{Recall}(C, \hat{R}) + \frac{\epsilon}{m} \geq \operatorname{Accuracy}(C, D)$

799 *Proof of Corollary 6.1.* By Corollary 4.1, $|\operatorname{Recall}(C, R) - \operatorname{Recall}(C, \hat{R})| \leq \frac{\epsilon}{m}$. This can be written
800 as $\operatorname{Recall}(C, \hat{R}) + \frac{\epsilon}{m} \geq \operatorname{Recall}(C, R)$. Using Corollary 2.1, $\operatorname{Recall}(C, \hat{R}) + \frac{\epsilon}{m} \geq \operatorname{Recall}(C, \hat{R}) + \frac{\epsilon}{m}$
801 $\geq \operatorname{Recall}(C, R) \geq \operatorname{Accuracy}(C, D)$. \square

Table 5: Notable Private Malware Reference Datasets

Name	Samples	Families	Platform	Collection Period
Malsign	142,513	?	Windows	2012 - 2014
MaLabel	115,157	?	Windows	Apr 2015 or earlier
MtNet	1.3 million	98	Windows	Jun 2016 or earlier

Table 6: Notable Public Malware Reference Datasets

Name	Samples	Families	Platform	Collection Period
VX Heavens	271,092	137	Windows	?
Malheur	3,133	24	Windows	2006 - 2009
MalGenome	1,260	49	Android	Aug 2010 - Oct 2011
Drebin	5,560	179	Android	Aug 2010 - Oct 2012
Malicia	11,363	55	Windows	Mar 2012 - Mar 2013
Kaggle	10,868	9	Windows	Feb 2015 or earlier