

# Title

Sanaa Siddiqui, Diganta Mukhopadhyay, Afzal Mohammed, Hrishikesh Karmarkar, Kumar Madhukar

## ACM Reference Format:

Sanaa Siddiqui, Diganta Mukhopadhyay, Afzal Mohammed, Hrishikesh Karmarkar, Kumar Madhukar . 2024. Title. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, ?? pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

TODO: Add affiliations

## 1 Introduction

TODO: Re-work wrt new experiments section

Advances in Deep Neural Networks (DNNs) have enabled the scalable solution of several previously intractable problems like image recognition. TODO: cite Due to this, DNNs have been increasingly assumed a central role across various domains. These include several safety critical domains like healthcare [b1], where they contribute significantly to medical diagnostics and predictive analysis [b2], and autonomous vehicles, where DNNs serve as the backbone for sophisticated perception systems, supporting tasks such as object recognition and decision-making [b3].

However, DNNs are well known to be vulnerable to adversarial attacks TODO: cite and can produce unexpected behaviours which in safety critical domains can lead to unsafe situations. Therefore, there is a critical need to guarantee the precise functionality and safety of these DNNs that are deployed in safety-critical environments, to build trust on the systems that utilize these DNNs.

A number of verification techniques have been proposed to validate the reliability of DNNs in safety-critical settings [b4, b5], including those based on abstract interpretation, branch and bound, etc. TODO: cite However, since DNN verification is NP-Hard TODO: cite, verification techniques for DNNs often face scalability issues. Our work is based on counterexample guided abstraction refinement [cegar-nn] and aims to alleviate this issue by reducing the verification problem for a large *concrete* DNN  $\mathcal{N}$  into an over-approximate query involving a smaller *abstract* DNN  $\mathcal{N}'$ . This is done by *merging* groups of neurons belonging to the same *merge*

group in  $\mathcal{N}$  into single neurons in  $\mathcal{N}'$ . TODO: appeal to later sections via ref. Then, this smaller query may be dispatched using any existing technique, thus extending the scalability of that technique.

While counterexample guided abstraction refinement based on structural abstraction has been proposed [cegar-nn] TODO: cite others, in several cases these techniques are forced to refine all the way to the original network, resulting in a query that is as complicated as the original verification problem we started with. We posit that this is because the refinement process does not identify which merges should be prioritised based on the semantics of the network, instead removing a single node from a merge group, leaving potentially equally over-approximate merges in the network. This leads to situations where, after several refinement steps,  $\mathcal{N}'$  has a large number of singleton merge groups, leading to a large yet over-approximate abstract network. TODO: refer to later sections with eg demoing this.

In this work we utilise measures of similarities between neurons explored in TODO: cite, check Jan Kretinsky's and others works to build a partial order of possible merges in each layer capturing the priority each merge should get during the refinement process. This allows us to use the information from a spurious counterexample to cut a merge group in a manner that eliminates the merge operation that has been estimated to contribute most to the over-approximation. This avoids the proliferation of singleton merge groups in  $\mathcal{N}'$ , and leads to a smaller and a tighter abstraction.

We describe a way to efficiently implement this partial order based refinement procedure to be able to iterate through potential refinements with minimal computational overheads. We do this by creating tree-based data structures that allows us to pre-compute a significant portion of the relevant information at the beginning of the refinement loop. Then, refinement operations can be performed via very quick index manipulations. This allows us to iterate through and screen several possible refinements fast before committing to a solver call.

## 2 Background

An example of Elboher's Refinement where we get a lot of singletons

## 3 Our Method

In same example, we immediately arrive at a better picture. This picture is equivalent to refining and merging back in the Elboher's setting.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference'17, July 2017, Washington, DC, USA

© 2024 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM. . \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 4 Methodology

Our methodology involves two broad steps:

1. Finding a tree structure that represents the order in which neurons should be merged.
2. Using that structure to guide a CEGAR approach in order to help reduce number of refinement steps.

tree\_merges.pdf

Order\_of\_Merge.pdf

### 4.1 Tree and Merges:

To establish the merging order of neurons, we create a tree structure wherein leaf nodes represent the original neurons, and non-leaf nodes represent merge groups. The construction of the tree follows a bottom-up approach, prioritizing the merging of similar neurons and delaying the merging of dissimilar ones. Similarity is ascertained through observation vectors. Consider simulating the network with distinct input values and observing the fluctuating values of a particular neuron corresponding to those inputs. The collection of these observed values constitutes an observation vector.

**Figure 1.** Order Of Merging

In Figure ??, when we input values  $\{1, 2, 3\}$  to the neuron  $n_0^0$  at different time instances, we obtain corresponding values  $\{1.5, 3, 4.5\}$ , forming the observation vector at neuron  $(n_1^1, Inc)$ . For instance, considering three neurons  $n_i^w$ ,  $n_j^w$ , and  $n_k^w$  with observation vectors  $v(n_i^w)$ ,  $v(n_j^w)$ , and  $v(n_k^w)$ , the merging sequence adheres to the following conditions:

If  $\|v(n_i^w) - v(n_j^w)\|_2 \leq \|v(n_i^w) - v(n_k^w)\|_2$  and  $\|v(n_i^w) - v(n_j^w)\|_2 \leq \|v(n_j^w) - v(n_k^w)\|_2$ , then  $n_i^w$  and  $n_j^w$  are initially merged into a representative neuron  $\alpha$ , followed by the

merging of  $\alpha$  and  $n_k$ . Here  $\|v(n_i^w) - v(n_j^w)\|_2$  computes the “Euclidean Distance” between the observation vectors  $v(n_i^w)$  and  $v(n_j^w)$ .

In Figure ??, for the first layer, the initial merge would involve the neuron  $(n_0^1, Inc)$  with  $(n_1^1, Inc)$ , forming *Merge 1*, then we would combine  $(n_2^1, Inc)$  and  $(n_3^1, Inc)$ , forming *Merge 2*. Finally we combine *Merge 1* and *Merge 2* into the root node.

As we progress up the tree, the degree of over-approximation rises. This is due to the increasing difference between observation vectors as we ascend. Therefore, the sub-trees closer to the root are indicative of coarser merges, whereas the ones farther from the root represent finer merges.

## 4.2 Using Counter-Examples to make cuts in the Tree

We are guided by this tree as a prospective refinement method. Starting with the entire tree where everything is merged. When the solver detects a spurious counterexample ‘ $\beta$ ’, we leverage it to refine the network. This process commences by identifying the “culprit neuron  $\gamma$ ” selected for refinement. A “culprit neuron” in a merge group is selected on the basis of how much the neuron contributed to the output. If change in output of neuron changes the value of the output neuron significantly then that neuron is a good candidate for “culprit neuron”.

Following this, we reverse all merges dependent on the culprit neuron  $\gamma$ . Therefore, refinement essentially involves finding a cut-point in the tree, precisely where all merges dependent on the culprit neuron  $\gamma$  are undone. Each cut produces a set of trees, the merge groups then consist of neurons in the leaf nodes of the these trees. Therefore finding new merge groups for refinement is therefore just finding a cuts in the tree.

Consider Figure ??, illustrating the merging sequence of neurons  $(n_0^1, Inc)$ ,  $(n_1^1, Inc)$ ,  $(n_2^1, Inc)$ , and  $(n_3^1, Inc)$ . If, for instance, the neuron  $(n_3^1, Inc)$  is identified as the problematic neuron based on a counter-example, we will reverse all the merges dependent on the  $(n_3^1, Inc)$  neuron, including *Merge 2* and the *Root Node* merge. Consequently, after implementing this reversal indicated in Figure ??, our refinement phase will yield three distinct merge groups. The first merge group comprises two neurons, namely  $(n_0^1, Inc)$  and  $(n_1^1, Inc)$ . The second merge group and the third merge have single neurons  $(n_2^1, Inc)$  and  $(n_3^1, Inc)$ , respectively.

## 4.3 Culprit Neurons

A neuron, denoted as  $\gamma$ , is designated as a culprit neuron within a specific layer when absolute value of the product of the difference between  $(v_{Rep(\gamma)})$  and  $v_\gamma$  and the effective weight is maximized.

$$\|(v_{Rep(\gamma)} - v_\gamma)\|_2 \cdot |(effective\_weight)|$$

In this context, *Rep* signifies the representative neuron for neuron  $\gamma$ ,  $v_\gamma$  represents the value of the neuron  $\gamma$  at

\_before and after cut.pdf

**Figure 2.** Trees and Cuts

counter-example  $\beta$  and *effective\_weight* represents the how much does the value of output neuron changes with respect to change in the value of the neuron under consideration, essentially corresponding to the “gradient” at that particular counter-example “ $\beta$ ”.

[1]  $\gamma = \arg \max_i \|v_{Rep(i)} - v_i\|_2 \cdot |effective\_weight|$  Find a sequence of nodes,  $t_1, t_2, t_3, \dots, t_k$  representing a path from  $t_1 = \text{root}$  to  $t_k = \gamma$ . Remove the nodes  $t_1, t_2, \dots, t_{k-1}$  denoting the merges dependent on  $\gamma$  through this path, leading to our connected tree being split into a collection of disconnected sub-trees. New merge groups are the leaf nodes in our disconnected graph. **Output** New Merge Groups

#### 4.4 Optimality of the Trees

Our objective is to determine the most efficient order for merging neurons, minimizing the introduction of over-approximation at each step. This approach aims to avoid creating networks with excessive over-approximation, which could lead to the generation of spurious counter-examples in response to queries. Opting not to mitigate over-approximation at each step would result in an increased number of refinement steps. This essentially entails making additional solver calls, incurring significant costs to eliminate the spurious counter-examples.

Nevertheless, during the initial merging process (until saturation is reached), the root node " $\rho$ " will exhibit the same level of over-approximation across all conceivable merging scenarios—for all possible tree sequences. Nevertheless, when we descend one level down the tree to explore the children nodes of our original root node  $\rho$  for the purpose of identifying a cut for refinement, we discover varying levels of over-approximation manifesting in the root nodes of the resultant sub-trees. These differences are a result of the different merging scenarios pursued to construct those individual trees.

[H] [1] Initialize every simulated distance vector as a singleton cluster. Initialize  $C = \{v_1, v_2, v_3, \dots\}$  as the set of singleton clusters. Initialize a Binary Tree  $T$  with leaves as  $\{(n_1), (n_2), (n_3), \dots\}$  corresponding to  $\{v_1, v_2, v_3, \dots\}$ . Initialize  $V$  as a set of visited nodes, empty at first.

MergeFunction  $u, v$  All nodes are classified as **Inc**

$\max(u, v)$

$\min(u, v)$

$|C| > 1$   $v_j, v_j = \arg \min_{a, b \in C} \|a - b\|_2$  Set  $w = \text{MergeFunction}(v_i, v_j)$

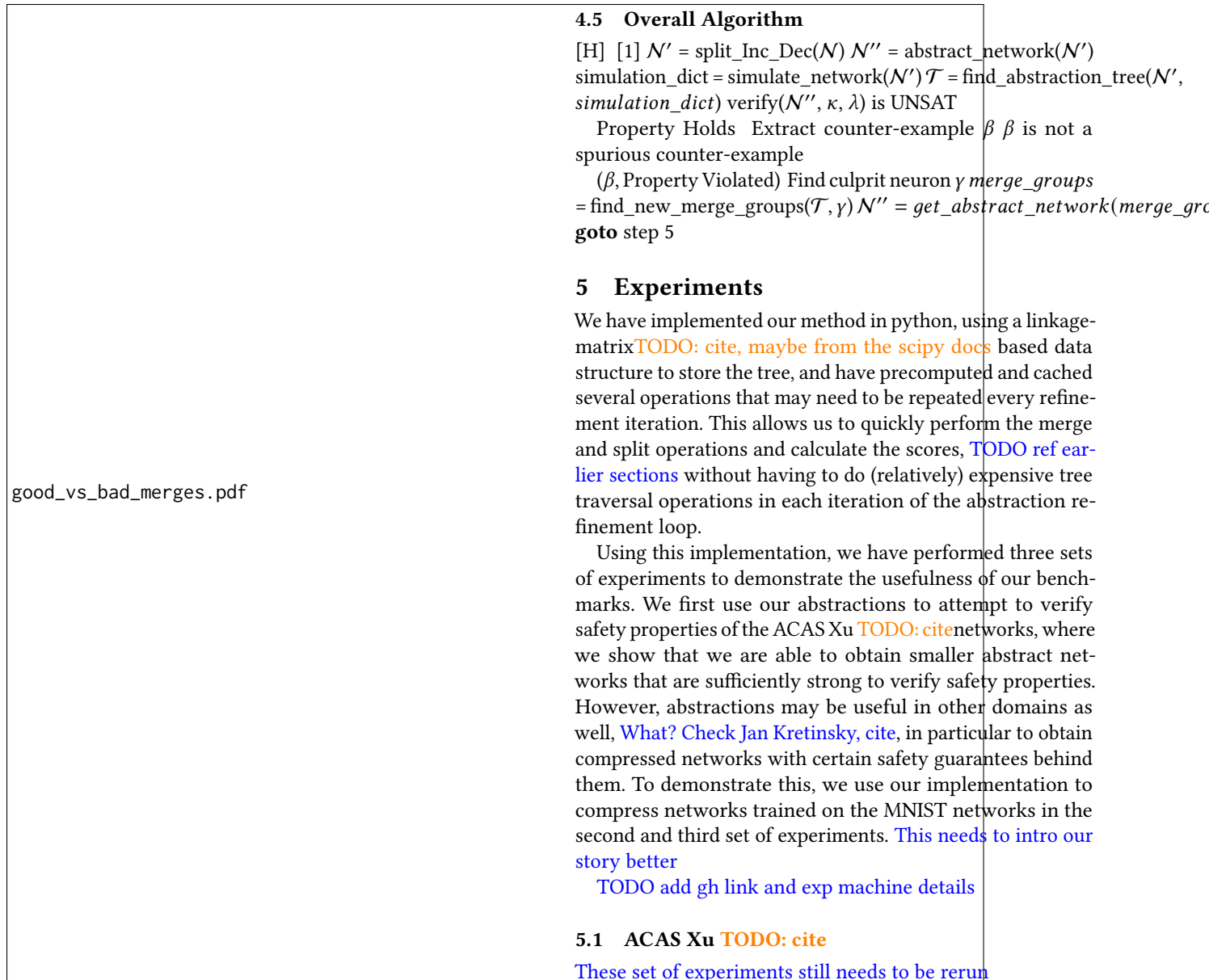
Let nodes from  $T$  not in  $V$  corresponding to  $v_i, v_j$  be  $m_i$  and  $m_j$  Remove  $v_i, v_j$  from  $C$  and add  $w$  to  $C$ . Make  $(m_i \cup m_j)$  the parent of  $(m_i)$  and  $(m_j)$  in tree  $T$  Add  $m_i$  and  $m_j$  to  $V$ .

While the optimal tree, representing the optimal merging sequence, can aid in the refinement process by guiding the reversal of merges, finding such an optimal tree poses is extremely challenging. Even when dealing with only 'n' Increment (Inc) neurons that have been merged to saturation, the total number of possible trees is given by  $(2n - 3)!!$ , making the task of determining the truly optimal tree from these options extremely challenging.

Since finding this ideal tree is a challenging task, we employ hierarchical clustering (Algorithm ??) as an approach to approximate and derive such a tree. Initially, we simulate our network using a set of ' $k$ ' inputs. Subsequently, we employ cluster analysis on these ' $k$ ' points to construct a hierarchical arrangement of clusters. This process initiates with data points corresponding to simulated values (observation values in the observation vector) of a neuron forming their own cluster. The clusters are then systematically combined based on their similarity, thereby generating a hierarchy of clusters.

The choice of similarity measure is the "*distance metric*" between clusters. We have used "*Euclidean Distance*" as our distance metric. Given that the data points to perform this hierarchical clustering originate from the values of the simulated neurons, this hierarchical clustering effectively reflects the methodology we employ to merge the neurons.

For example, in Figure ??, we conducted a simulation of our network on three data points. Subsequently, we examined the observation vectors corresponding to these points. Utilizing the hierarchical clustering algorithm, the initial selection for merging will involve  $(n_0^1, Inc)$  and  $(n_1^1, Inc)$  because of the fact that their Euclidean distance is minimum. This forms *Merge 1* in Figure ??. The observation vector for *Merge 1* ( $v_{Merge1}$ ) is the max of the  $v((n_0^1, Inc))$  and  $v((n_1^1, Inc))$  which is  $\{1.5, 3, 4.5\}$ . For decrement nodes the observation vector would be minimum of the observation vector of the corresponding decrement nodes. The next merging step involved selecting  $v((n_2^1, Inc))$  and  $v((n_3^1, Inc))$  and merging these two neurons, representing *Merge 2* in Figure ??. The observation vector for *Merge 2* is now  $\{4, 8, 12\}$ . Ultimately, the *Merge1* merge group is merged with the *Merge 2* merge group to create the Root Node in our network.



**Figure 3.** Ways of Merging

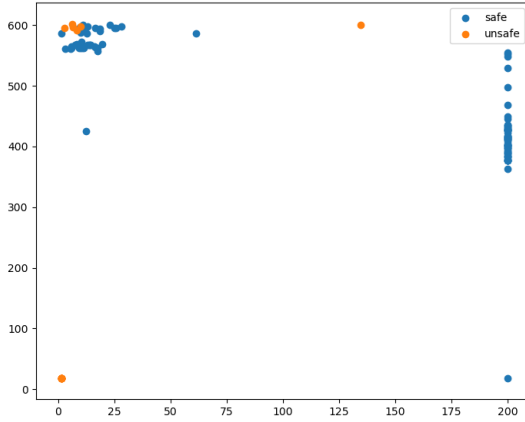
This approach of merging neurons based on similarity proves advantageous as it helps in reducing number of refinement steps. For instance, consider the task of checking whether  $\forall v_0^0 \in [0, 1]$  implies  $v_0^2 < 10$ . If we had the neuron  $(n_3^1, Inc)$  as the culprit neuron and then if we follow the second merging approach depicted in Figure ??, then we would have been compelled to reverse both *Merge 1* and *Merge 2*. However, by employing the first merging approach, undoing only the *Merge 2* becomes sufficient, resulting in a reduction in the number of refinement steps required.

In these set of experiments, we have taken the ACAS Xu [TODO: cite](#) set of networks [TODO: add detailed table](#) and properties and attempt to verify them using a CEGAR [TODO: cite](#) approach. We use our technique to generate the abstraction, and attempt to verify the property on the abstract network using an existing neural network verification solver. If the solver returns a counterexample, and if that counterexample is not spurious, we use that counterexample as a reference point for score calculation [TODO: Ref prev sec](#) to select the culprit neuron. [TODO: Ref prev section](#) With the culprit neuron selected, we then refine the network.

The results are summarised in table ??. The solver we used was  $\alpha\beta - CROWN$  [TODO: cite](#). Note, avg time and network size here doesn't include timeouts. Also, for reduction rate, the original network size is taken to be 600, since the re-merging optimization is not there in this data. With that

Benchmarks	% Verified	Avg Reduction	Avg Time
safe	40.404	4.4875	14.198
unsafe	100	55.6071	12.7916

**Table 1.** Summary of ACAS Xu [TODO: cite](#) Results [TODO: Fix formatting](#)



**Figure 4.** Scatter plot of abstract network sizes vs time in seconds for ACAS Xu [TODO: cite](#)

optimization added, we should compare against 300. Also, the full verification rate is 100% because its not a full set of experiments, this will potentially be a lot lower.

We notice that for the unsafe cases, we are able to refute the properties via a counterexample obtained on a significantly smaller network. But what about the safe cases, what do we say?. In fact, for most networks, we are able to refute based on a counterexample from the fully abstracted network of size 18. This is not there in this small summary table, but is clear from the full data table. Maybe link to that somewhere in the appendix? Note that since we use only two classes in our fully abstracted network as opposed to 4 in [cegar-nn], the size of this network is significantly smaller. In other networks, we are able to get a counterexample on a network with size very close to the original network.

For the safe cases, we see that the size of the abstracted networks are large for instances that did not time out. This is in line with what was observed in [cegar-nn] Can we claim this? It's not immediately clear from their paper, maybe we need to show our recreation of their experiments? Also, this split is not clear from the table. For other instances, the timeout happened due to the solver timing out on the abstract query. While the sizes of the networks on which this to happened are small, does that not say something negative about our abstraction?

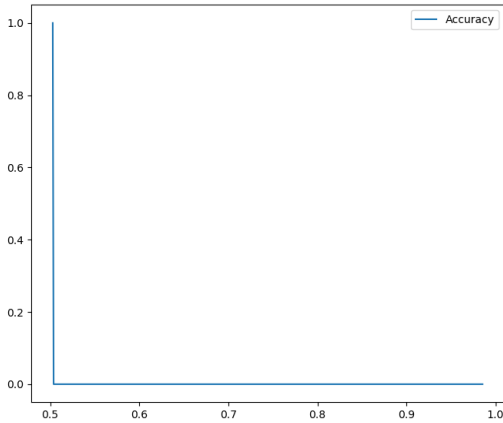
Figure ?? shows the network sizes and times of the ACAS Xu [TODO: cite](#) networks as a scatter plot. This includes to and non-to. Note that the time taken here is dominated by the solver call time Should we measure and report this precisely? We see that the time taken by the solver call on these networks is not at all correlated with the network sizes. For instances, for the safe cases, there are several networks of sizes ranging from 350 to 600 ReLU nodes all of which time out on the solver call. Similarly, for the unsafe networks, several 600 size networks take less than 25 seconds. In fact, we found that the time taken by the solver call on a particular network is more dependent on the solver (and configuration of the solver) used than on the size of the network. Due to this, in the following experiments, we measure the accuracy of the networks obtained via abstraction instead of making solver calls to determine the effectiveness of the abstraction. Is this okay? Can it be instead argued that the variance is due to the networks themselves, so the abstraction is not necessarily useful? Maybe trying with other configs and showing the variance?

## 5.2 Robustness on VNNComp MNIST [TODO: cite](#)

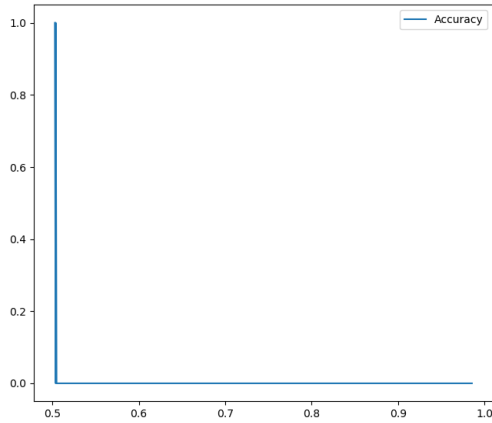
What about doing these experiments on ACAS? Do pgd then samples

In this section, we start with our abstraction and progressively refine it using our refinement technique, measuring the number of spurious counterexamples at each step. To perform the refinement at each step, we try four methods to perform the culprit neuron selection. The first method 'random' simply chooses a culprit at random, and serves as a baseline. For the other three methods, we choose a number of reference points to calculate scores, and choose the culprit neuron with the highest score. For 'samples', the reference points are random generated input points that satisfy the precondition. For 'pgd', the samples are generated via a PGD [TODO: cite](#) attack on the abstract network at each step. For 'samples-pgd', we do some initial refinement steps using the randomly generated samples, when all the samples have been exhausted, we then perform 'pgd'. We use the implementation of 'pgd' inside  $\alpha\beta$ -CROWN [TODO: cite](#) for our experiments.

Figure ?? ?? shows a typical plot for these experiments. Here accuracy measures the number of true counterexamples within a random sample drawn from a region satisfying the pre-condition. Should I plot and talk in terms of spurious counterexample rate instead? We see that as we refine the network and the reduction rate reduces (moving right to left along the graph), the accuracy remains close to zero. At some point however, the accuracy jumps up sharply and becomes close to 100%. This indicates that there are a few critical refinement steps that remove almost all spurious counterexamples. This may be due the fact that with epsilon robustness properties that are being explored here, the pre-condition region may be very small, so a single refinement step may



**Figure 5.** Accuracy vs Reduction Rate plot for VNNComp MNIST of size  $2 \times 256$  with  $\epsilon$ -Robustness property. Refinement done via the 'samples' method.



**Figure 6.** Accuracy vs Reduction Rate plot for VNNComp MNIST of size  $2 \times 256$  with  $\epsilon$ -Robustness property. Refinement done via the 'pgd' method.

change the behavior of the network on that small region in a way that eliminates all possible spurious counterexamples from that region. Is this expl okay? BaB and abstract interp based methods like  $\alpha\beta - CROWN$  are successful on such props because of this, discuss this? Also, other graphs are present. For eg, some random sampling graphs show peaks in the middle (2..random), and some pgd graphs are flat lines (4,6..pgd). Discuss these?

Table ?? summarizes these results for multiple VNNComp MNIST networks and multiple refinement methods. Accuracy here refers to the final accuracy of the best

Net Size	Cex Method	Reduction	Accuracy	No. Steps
$2 \times 256$	random	49.1%	100%	809
$4 \times 256$	random	49.6%	100%	1778
$6 \times 256$	random	49.7%	100%	2602
$2 \times 256$	samples	50.3%	100%	27
$4 \times 256$	samples	49.9%	100%	648
$6 \times 256$	samples	47.9%	100%	1199
$2 \times 256$	pgd	50.4%	100%	28
$4 \times 256$	pgd	88.6%	0%	14
$6 \times 256$	pgd	83.3%	0%	82
$2 \times 256$	samples-pgd	50.3%	100%	27
$4 \times 256$	samples-pgd	49.9%	100%	648
$6 \times 256$	samples-pgd	47.9%	100%	1209

**Table 2.** Summary of VNNComp MNIST Results on a single robustness property

network found when the refinement process stops due to lack of counterexamples. Firstly, we find that 'pgd' performs poorly on the 4 and 6 layer networks, where while the number of steps are relatively low and the reduction rate is high, the accuracy is very low. This is because the  $\alpha\beta - CROWN$  implementation of PGD-attack fails to find a spurious counterexample early on in the refinement process for these instances, leading to early termination without ever hitting the accuracy jump.

Apart from these exceptions, the other instances and methods seem to perform comparably on the same network in terms of reduction rate. This is evidence for the fact that, within the limited search space of abstractions defined by our tree, there exists a strong enough abstract network with a good reduction rate, and that it is possible to find this network via a CEGAR-like approach.

Further, we see significant difference in the number of steps taken by 'random' and the other methods in finding the abstract network. This indicates that our heuristics for finding a culprit neuron to base the refinement on indeed shortens the search for the refined network. We also observe that 'samples', 'samples-pgd' and for the 2 layer network 'pgd' take similar number of refinement steps, indicating that each method may be taking close to the optimal set of refinement steps. So, although the 'pgd' method incurs a significant time overhead it does not necessarily provide any benefits. Last two paragraphs okay? Refer to this in intro of experiments section? Also, pgd does not produce improvement (unlike in cleverest). This is prob because global view, and cegarette is about doing an eager attack before solver call. Discuss this? So two potential angles on this, what Ive written vs global view angle.

### 5.3 Compression with Guarantees on Critical Classes

In several safety-critical applications of DNNs classifiers, there are certain 'critical' classes for which a false negative

classification is far more dangerous than a false positive one. For example, for medical diagnosis and collision detection, a false negative is far more dangerous than a false positive.