
ARACHNE: SEARCH BASED REPAIR OF DEEP NEURAL NETWORKS

A PREPRINT

Jeongju Sohn*

University of Luxembourg
Luxembourg city
Luxembourg
jeongju.sohn@uni.lu

Sungmin Kang*

KAIST
Daejeon
Republic of Korea
sungmin.kang@kaist.ac.kr

Shin Yoo

KAIST
Daejeon
Republic of Korea
shin.yoo@kaist.ac.kr

ABSTRACT

The rapid and widespread adoption of Deep Neural Networks (DNNs) has called for ways to test their behaviour, and many testing approaches have successfully revealed misbehaviour of DNNs. However, it is relatively unclear what one can do to correct such behaviour after revelation, as retraining involves costly data collection and does not guarantee to fix the underlying issue. This paper introduces Arachne, a novel program repair technique for DNNs, which directly repairs DNNs using their input-output pairs as a specification. Arachne localises neural weights on which it can generate effective patches and uses Differential Evolution to optimise the localised weights and correct the misbehaviour. An empirical study using different benchmarks shows that Arachne can fix specific misclassifications of a DNN without reducing general accuracy significantly. On average, patches generated by Arachne generalise to 61.3% of unseen misbehaviour, whereas those by a state-of-the-art DNN repair technique generalise only to 10.2% and sometimes to none while taking tens of times more than Arachne. We also show that Arachne can address fairness issues by debiasing a gender classification model. Finally, we successfully apply Arachne to a text sentiment model to show that it generalises beyond Convolutional Neural Networks.

1 Introduction

Deep Neural Networks (DNNs) are rapidly being adopted in many application areas [1], ranging from image recognition [2, 3], speech recognition [4], machine translation [5, 6], to safety-critical domains such as autonomous driving [7, 8] and medical imaging [9]. As the application areas expand, there has been a growing concern that these DNNs should be *tested*, both in isolation and as a part of larger systems, to ensure dependable performance. The need for testing resulted in two major classes of techniques. First, to evaluate sets of test inputs, many test adequacy criteria have recently been proposed [10, 11, 12]. Second, ways to synthesise new inputs by applying small perturbations to given inputs (such as emulation of different lighting or weather conditions) have been introduced [13, 14, 15]. Newly synthesised inputs not only increase input diversity, but also can reveal unexpected behaviour of DNNs under test.

Compared to traditional software developed by human engineers, however, the stages after the detection of unexpected behaviour remain relatively unexplored for DNNs. This is due to the major difference between the way DNNs and code are developed: one is trained based on data, while the other is written by human engineers based on specifications. Consequently, existing efforts to *repair* the unexpected behaviour tend to heavily rely on retraining [16]. However, the use of retraining as a means to repair a DNN has a couple of weaknesses. First, since retraining uses the overall accuracy of a DNN to assess the quality of a repair, it may not remove the unexpected behaviour even if retraining increases overall accuracy. Apricot, a state-of-the-art DNN repair technique [17], fixes DNNs iteratively by incorporating training into its repairing process. As Apricot also aims to increase general accuracy, it is not geared towards removing unexpected behaviour and may have difficulty removing a specific misclassification pair. Second,

*First co-authors with equal contribution

systematic retraining can be computationally expensive. In the case of Apricot, there is a training loop within the repairing process, which causes Apricot to take multiple hours repairing DNNs.

This paper introduces Arachne, a search-based automated program repair (APR) technique for DNN classifiers. Compared to existing techniques, Arachne is more focused on fixing specific misclassifications, e.g. erroneously predicting an input of class a as class b , rather than overall accuracy. Arachne directly searches the space of neural weights, guided by a novel fitness function inspired by Generate and Validate APR techniques [18, 19, 20]. Arachne resembles APR techniques for traditional code in many aspects: it attempts to localise components relevant to the given misclassification, and uses both positive (i.e., inputs that are correctly classified) and negative inputs (i.e., inputs that are misclassified) to retain correct behaviour and to generate a patch, respectively. As Arachne directly changes the values of neural weights, its internal representation of a patch is a vector of real numbers. Arachne thus adopts Differential Evolution (DE) [21], a meta-heuristic optimisation algorithm shown to be effective in continuous search spaces [21, 22], as its search algorithm. During each fitness evaluation, Arachne updates the localised neural weights of the DNN under repair with values with the DE candidate solutions. Subsequently, it evaluates the fitness of solutions based on the classification results of inputs.

We empirically evaluate Arachne on four image classification benchmarks (i.e., Fashion-MNIST [23], CIFAR-10 [24], German Traffic Sign Recognition Benchmark (GTSRB) [25] and Labelled Faces in the Wild (LFW) [26]) and one text classification benchmark (i.e., Twitter US Airline Sentiment dataset [27]). We use the first three benchmarks to show the feasibility of fault localisation and patch generation. We also use them to compare Arachne to Apricot: for CIFAR-10, we reuse the three CNN architectures employed to evaluate Apricot [17] for fair comparison and build an additional image classifier model with more training capacity than the one used to evaluate the feasibility individually for Fashion-MNIST and GTSRB. The LFW benchmark is an award-winning benchmark of human faces: we use it to construct a case study that shows how Arachne can fix a fairness issue in a gender classifier by rebalancing the model. For the Twitter benchmark, which contains various tweets and their underlying sentiments, we employ it to evaluate whether Arachne can address the models other than those for image classification.

The technical contributions of this paper are as follows:

- The paper introduces Arachne, a novel search-based repair technique for DNNs. Unlike existing approaches that retrain a DNN model using more inputs, Arachne aims to repair a pre-trained model by directly adjusting neural weights.
- We empirically evaluate Arachne against the state-of-the-art DNN repair technique, Apricot, using widely studied image classification benchmarks. Arachne can produce repairs that are more focused on the targeted misclassifications, while only minimally perturbing other behaviour, and operate at speeds tens of times faster than Apricot.
- We present a case study on how Arachne can be used to repair a DNN model suffering from a fairness issue. Arachne can successfully repair bias in the given gender classifier without requiring any additional data. Arachne improved the classification accuracy for female images from 86.8% to 88.8%.
- We conduct an additional study with a model that analyses the underlying sentiment of text in order to evaluate whether Arachne can be effective in different domains besides image classification. Arachne successfully decreases the prevalence of the most frequent errors, showing that Arachne can work with diverse models.

The remainder of this paper is organised as follows. Section 2 describes components of our DNN repair technique, Arachne. Section 3 sets out the research questions and describes experimental protocols. Section 4 outlines the set-up of the empirical evaluation, the results of which are discussed in Section 5. Section 6 presents threats to validity, Section 7 details the related work, and Section 8 concludes.

2 Arachne: Search Based Repair for DNNs

This section motivates the development of Arachne and describes its internal components.

2.1 Motivation

The standard approach towards getting rid of misbehaviour in DNNs is through further training, which typically involves a large volume of training data, curated carefully to avoid the misbehaviour. Unfortunately, preparing training data is known to be the major bottleneck to the practical application of machine learning due to its high cost [28, 29]. Furthermore, even with new curated training data, retraining does not always remove the observed misbehaviour, as retraining aims to improve overall accuracy instead.

We posit that there are situations that call for *targeted* improvements of DNN models with more assurance for the improved outcome. For example, certain misclassifications, such as mistaking a stop sign for a 60 kilometers per hour speed limit sign, may pose more risk than other mistakes about parking allowance. One may also argue that some misbehaviour that stem from inherent bias in the data should be fixed, even if it degrades overall accuracy, as in the case of gender bias in facial classification [30]. Presented with such misbehaviour that needs to be repaired urgently, retraining loses its appeal as a correction mechanism due to the cost of data curation. In such a scenario, a more direct approach that does not require additional data may be more useful.

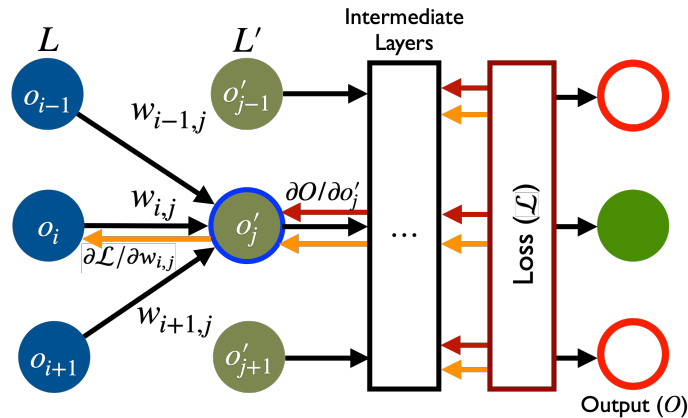
We emphasise that our aim with Arachne is not to replace well-designed learning processes as a means of improving general learning capability and overall model accuracy. Rather, we want to introduce and evaluate an alternative repair technique that can introduce a directed and focused improvement over a small set of unexpected behaviour, without requiring larger and better-curated training datasets. In this sense, we expect Arachne to complement other training-based techniques. Note that Arachne targets misbehaviour that can be repaired by neural weight manipulation alone: if a DNN model underperforms due to its inappropriate architecture, we do not expect Arachne to be effective at repairing it.

2.2 Overview

Arachne has two primary operations: localisation and patch generation. In the localisation phase, Arachne identifies a set of neural weights that are likely to be related to the observed misclassification revealed by a negative input set. The intuition is that changing the values of these neural weights is likely to affect the model behaviour in the desired direction.

As an effective repair should fix misbehaviour while minimally disrupting correct behaviour, Arachne leverages positive inputs (i.e., inputs that are processed correctly) in addition to negative inputs (i.e., inputs that reveal the misclassification) for both localisation and patch generation. The basic intuition behind fault localisation in software engineering is that faulty elements relate more to a program’s faulty behaviour (e.g., failing tests) and less to its correct behaviour (e.g., passing tests) [31]. Under the same intuition, Arachne localises neural weights that have more impact on negative inputs and less impact on positive inputs. In the patch generation phase, similarly to Generate and Validate (G&V) Automated Program Repair techniques such as GenProg [32] and Arja [19], the positive inputs are used by the fitness function of Arachne to retain the initially correct behaviour of the DNN under repair. Employing fitness values based on the inference results of both positive and negative inputs, Arachne uses Differential Evolution [21] to search for a set of neural weights that would correct the behaviour of the DNN under repair. Details of the localisation and patch generation phases are discussed in the following subsections.

Figure 1: Arachne fault localisation diagram. Black arrows represent forward pass, while red and orange arrows are backpropagation. Black and red arrows combined compute the forward impact, whereas the orange arrow denotes the gradient loss. The gradient loss and the forward impact combined evaluate the likelihood of $w_{i,j}$ being localised by Arachne.



2.3 Localisation Phase

Attempting to adjust all neural weights of a DNN would be costly, as even relatively simple DNNs consist of thousands of weight parameters. Hence, instead of targeting all weights of the DNN, Arachne incorporates a novel localisation

method, Bidirectional Localisation (BL), that aims to identify neural weights on which Arachne can more effectively repair. As depicted in Figure 1, BL considers the gradient loss and the forward impact of each neural weight; the former quantifies the responsibility that a neural weight has for the misclassification, and the latter measures the influence of the weight on the final classification outcome. Since a weight that has a high impact is not necessarily responsible for the misclassification, BL treats the forward impact and the gradient loss as two competing objectives to optimise.

Algorithm 1 presents the pseudo-code of the bidirectional localisation method. The algorithm calls **ComputeGradientLoss** and **ComputeForwardImpact** to compute the gradient loss (Line 4-5) and the forward impact (Line 7-8) for each neural weight. Computing the gradient loss of a neural weight is straightforward. For instance, for a neural weight $w_{i,j}$ that connects the i th neuron in layer L , o_i , and the j th neuron in the following layer L' , o'_j , in Figure 1, its gradient loss backpropagated from the final loss, \mathcal{L} , is simply calculated as $\frac{\partial \mathcal{L}}{\partial w_{i,j}}$.

Computing the forward impact of a neural weight on the final output (**ComputeForwardImpact**) is more complicated than computing the gradient loss (**ComputeGradientLoss**). For a neural weight to have a significant impact on the final classification outcome, the weight has to be highly influential on the activation value of the output neuron to which it is connected (1), and this output neuron should be able to influence the final classification result (2). BL estimates the influence of a neural weight on the activation of the output neuron by multiplying the given neural weight and the average activation value of the neuron connected to the weight in the previous layer; it further normalises this estimated influence by dividing it by the aggregated influence of the neural weights connected to the same output neuron. Hence, in Figure 1, the influence of $w_{i,j}$ on the activation of o'_j (1) is computed as $o_i w_{i,j}$, which is then normalised as $\frac{o_i w_{i,j}}{\sum_l^{L'} o_l w_{l,j}}$ by dividing by the aggregated impact of all the neural weights connected to o'_j ($\sum_l^{L'} o_l w_{l,j}$). For the influence of the output neuron connected to the given weight (2), BL computes the gradient of this output neuron to the final output; in this example of $w_{i,j}$, this gradient value is calculated as the gradient of o'_j to the final output O (i.e., $\partial O / \partial o'_j$). After quantifying these two types of influence, (1) and (2), BL multiplies them to compute the forward impact of a neural weight at last; thus, for $w_{i,j}$ in Figure 1, the forward impact is $\frac{o_i w_{i,j}}{\sum_l^{L'} o_l w_{l,j}} \frac{\partial O}{\partial o'_j}$.

During the repair, we want to avoid unintentional disruption of the initially correct model behaviour as much as possible. Consequently, BL calls **ComputeGradientLoss** and **ComputeForwardImpact** twice, once with the negative inputs, I_{neg} (Line 4,7) and once with the positive inputs, I_{pos} (Line 5,8); for I_{pos} , instead of using all of them, we randomly sample the same number of positive inputs as the negative ones since a fully trained model often have far more positive inputs than the negative ones (Line 2). As shown respectively in Line 6 and 8, the gradient loss and the forward impact obtained with the negative inputs are divided by the corresponding values computed with the positive inputs². This ratio form aims to select neural weights related more to the negative inputs and less to the positive inputs.

The result of localisation is the set of weights that constitutes the Pareto front for the combined gradient loss and the forward impact of the positive and negative inputs (Line 12): that is, the set of weights such that there is no other weight that has both a greater gradient loss (ratio) and greater forward impact (ratio) simultaneously.

2.4 Patch Generation

Arachne uses Differential Evolution (DE) to generate patches that repair the misclassification of a DNN image classifier. Here, we describe how DE is configured for repair.

2.4.1 Differential Evolution

DE is a population-based metaheuristic optimisation algorithm that is known for its convergence speed and accuracy in continuous optimisation [21, 33, 22, 34]. DE was initially designed with multidimensional real-valued functions in mind [35], making it a good fit with Arachne whose aim is to quickly optimise a number of neural weights guided by a fitness function. A patch generated by Arachne is a set of new neural weights and is represented as a vector of dimension $N = |W_t|$, the number of neural weights returned by the localisation phase.

DE starts with a pool of candidate solutions and guides the pool to achieve higher fitness over multiple generations. At generation t , let $X^t = (x_1^t, \dots, x_N^t)$ be a target vector, i.e., a target solution under inspection, and V , a trial vector that may replace X^t in the next generation. To generate V , DE randomly selects three vectors from its current population: X_1^t , X_2^t , and X_3^t , all unique and different from the target vector X^t . Each element of V , v_i , is generated by adding the difference between $x_{2,i}^t$ and $x_{3,i}^t$ to $x_{1,i}^t$. Equation (1) shows how a candidate trial vector is created. DE employs two parameters, F and CR , to control the mutation rate and the cross-over rate, respectively. The mutation rate parameter, F , which is randomly sampled from a given range at every generation, controls the impact of the difference

²To avoid division-by-zero, one is added to the denominator.

Algorithm 1: Bidirectional Localisation of Neural Weights

input : A DNN model to be repaired, \mathbf{M} , the set neural weights within the DNN, W , a set of inputs that reveal the fault, I_{neg} , a set of inputs that correctly classified, I_{pos} , a loss function, \mathcal{L} , and the number of neural weight candidates to choose based on gradient loss, N_g

output: a set of neural weights to target for repair, W_t

- 1 $pool \leftarrow []$;
- 2 $I_{pos} \leftarrow \text{RandomSample}(I_{pos}, |I_{neg}|)$;
- 3 **for** $weight$ in W **do**
- 4 $grad_loss_{neg} \leftarrow \text{ComputeGradientLoss}(weight, \mathbf{M}, I_{neg}, \mathcal{L})$;
- 5 $grad_loss_{pos} \leftarrow \text{ComputeGradientLoss}(weight, \mathbf{M}, I_{pos}, \mathcal{L})$;
- 6 $grad_loss \leftarrow \frac{grad_loss_{neg}}{1+grad_loss_{pos}}$
- 7 $fwd_imp_{neg} \leftarrow \text{ComputeForwardImpact}(weight, \mathbf{M}, I_{neg})$;
- 8 $fwd_imp_{pos} \leftarrow \text{ComputeForwardImpact}(weight, \mathbf{M}, I_{pos})$;
- 9 $fwd_imp \leftarrow \frac{fwd_imp_{neg}}{1+fwd_imp_{pos}}$
- 10 Add tuple $(weight, grad_loss, fwd_imp)$ to $pool$;
- 11 **end**
- 12 $W_t \leftarrow \text{ExtractParetoFront}(pool)$;
- 13 **return** W_t

as mutation: the higher F is, the more mutated V will be. As shown in Equation 2, V replaces X^t if and only if V is better or at least equal to X^t . This selection strategy ensures that the population of DE improves monotonically.

$$v_i = x_{1,i}^t + F(x_{2,i}^t - x_{3,i}^t), i \in [1, N] \quad (1)$$

$$\vec{X}^{t+1} = \begin{cases} \vec{V}, & \text{if } fitness(\vec{V}) \geq fitness(\vec{X}^t) \\ \vec{X}^t, & \text{otherwise} \end{cases} \quad (2)$$

Algorithm 2 describes the pseudo-code of DE. Crossover rate, CR , in Line 8 to 10, controls the proportion of vector elements to be mutated: a lower CR value implies that only a small number of elements are mutated together, which results in DE searching each parameter separately [22]. Note that DE also ensures that at least one element of a trial vector is mutated (Line 5 and 8). We use classical DE, defining the range of F and the value of CR in advance, instead of adjusting them dynamically. Configuration details of DE are in Section 4.2.

2.4.2 Initialisation

The initial population of DE must be initialised. Since Arachne attempts to search for a patch by making adjustments to neural weights using DE, it may benefit from searching in the space near the original weights. Consequently, Arachne initialises the parameters of each candidate solution individually: the initial value of each parameter (i.e., a localised neural weight) is sampled from a Gaussian distribution, whose mean and standard deviation are computed using all the weights in the same layer with the parameter.

2.4.3 Fitness Function

Arachne uses a fitness function to evaluate each candidate solution and guide the search in DE. As in other G&V techniques, the fitness function of Arachne consists of two main parts: fixing misclassifications and retaining correct classifications. The fitness function for a candidate solution X is defined as follows:

$$score(i, X) = \begin{cases} 1, & \text{if } label_X(i) = label_{GT}(i) \\ \frac{1}{Loss(i, X) + 1}, & \text{otherwise} \end{cases} \quad (3)$$

$$fitness(X) = \sum_{i_p \in I_{pos}} score(i_p, X) + \alpha \sum_{i_n \in I_{neg}} score(i_n, X) \quad (4)$$

Algorithm 2: Differential Evolution

input : a mutation rate, F , a cross-over rate, CR , the size of population NP , the number of generations, Gen , and a fitness function, **fitness**

output: the best solution that has been found so far, $best$

```

1  $pop \leftarrow \text{Initialise}(NP)$ 
2 for  $gen$  in  $[1, Gen]$  do
3   for each  $X$  in  $pop$  do
4      $X_1, X_2, X_3 = \text{SelectRandom}(pop) // X_1 \neq X_2 \neq X_3 \neq X$ 
5      $j \leftarrow$  randomly selected between  $[1, N]$ 
6     for  $i$  in  $[1, N]$  do
7        $r_i \leftarrow \text{SelectUniform}(0, 1)$ 
8       if  $i = j$  or  $r_i \leq CR$  then
9          $v_i \leftarrow x_{1,i} + F(x_{2,i} + x_{3,i})$ 
10      end
11     end
12     if  $\text{fitness}(V) \geq \text{fitness}(X)$  then
13        $X \leftarrow V$ 
14     end
15   end
16 end
17  $best \leftarrow \text{SelectBest}(pop)$ 
18 return  $best$ 

```

For each input i , Equation (3) assigns a score based on whether it is correctly classified or not when a candidate solution X is applied. Here, $label_X(i)$ denotes the predicted label of an input i by a model patched with X , and $label_{GT}(i)$ is the ground truth label of the i . The score is inversely correlated to the loss when the input is incorrectly classified, to guide the search towards correct classification. However, once the input is correctly classified, the score is simply 1.0. This is to avoid overfitting, i.e., reducing the loss for correct images at the cost of ignoring classification results for misclassified images. The set I_{neg} contains inputs that reveal the misbehaviour, whereas the set I_{pos} contains inputs that are initially processed correctly. The overall fitness is the sum of scores of both sets of inputs; unlike the localisation, Arachne uses the entire set of I_{pos} to compute the fitness during the repair.

We use the hyperparameter α to balance the two components. A smaller α means more emphasis on preserving the correct behaviour, whereas a larger α means more emphasis on correcting the misbehaviour (shown by I_{neg}). Section 4.3 describes the details of how input sets, I_{neg} and I_{pos} , are composed, as well as how the default value for α is configured.

3 Experimental Protocols

This section outlines the experimental protocols for the empirical study and presents the research questions.

3.1 Faults for DNNs

This paper uses the term *repair* to refer to the process of correcting the misbehaviour of a DNN. The terminology for the cause of the misbehaviour can, however, lead to a philosophical question: is it possible to state that a DNN contains a *fault* when it misbehaves? While certain types of faults in DNN models are explicitly committed just as in code (e.g., an incorrect choice of layer [36]), the kind of DNN misbehaviour that we target is actually anticipated due to the nature of machine learning. This difference may have a nontrivial impact on future work on DNN testing and repair. Note that the use of the term fault and faulty input in this paper is for the sake of convenience, and is not intended to answer the question about the nature of faults in DNNs.

One practical issue with the type of “faults” we target is that it may not be easy to curate, document, and create a benchmark of them, as they are model specific and not explicitly *committed*. To circumvent the lack of fault benchmark, we evaluate Arachne with two classes of faults: those artificially injected by weight perturbation, for the evaluation of the localisation phase, and those naturally emerging after full training (i.e., inputs outside the training data that are *not* malicious attacks generated intentionally, but nonetheless cause misbehaviour of the given model), for the evaluation of the patch phase. We solely use the artificial faults to evaluate the localization effectiveness of Arachne because this

assessment requires the location of faults to be explicitly known. This is in contrast to patch effectiveness which is evaluated on naturally emerging faults; in this case, the exact fault location is not necessary for evaluation.

3.2 Research Questions

We investigate the following six research questions to evaluate the effectiveness of Arachne. The artificial faults are only used for the evaluation of RQ1, which evaluates the localisation effectiveness of Arachne; meanwhile, RQs 2 to 6 use naturally emerging faults to assess the patch generation phase of Arachne from diverse aspects.

3.2.1 RQ1. Localisation Effectiveness

How effective is the fault localisation of Arachne? To answer RQ1, we inject artificial faults into DNN models by weight perturbation. For this, we first randomly select a neural weight among all layers whose gradient to the final output is above average. Here, we target only one neural weight to ensure that the observed change in model behaviour is due to the perturbation of this weight. We then mutate the selected neural weight by adding noise from the standard normal distribution, $\mathcal{N}(0, 1)$, until it changes the behaviour of at least 0.1% of *all test inputs*; the change of behaviour includes both cases of misclassifying initial correct inputs and the opposite. We set this lower limit of changed model behaviour empirically. Since the impact of perturbing a single weight on the final classification is rather limited, the obtained limit is set at a low value, 0.1%. We further record the actual proportion of changed behaviour, Change Ratio, to inspect the correlation between this value and the effectiveness of localisation. After inserting a fault into the model through the mutation, we examine if our localisation method can identify this mutated neural weight.

We compare our method against two baselines: Gradient Loss (GL) based selection, a variation of our fault localisation method that uses only the gradient loss values of neural weights, as well as Random Selection (RS), where a random ordering is assigned to neural weights. We repeat the experiment 30 times to show that our fault localisation method works consistently.

3.2.2 RQ2. Patch Feasibility

Can Arachne repair misbehaviour of a DNN model? To show that Arachne can repair a DNN model by directly manipulating the neural weights, we randomly choose ten percent of the images that reveal misbehaviour of a DNN image classifier in *the test inputs* (i.e., those unseen during the training) and investigate whether Arachne can repair them. As the DE algorithm used by Arachne is inherently stochastic, we repeat this process 30 times. This repetition includes the uniform random sampling of misbehaviour, where each sample works as a different set of misbehaviour that a model can have. Hence, we posit that any potential bias in the selection of neurons will be minimised by this repetition. We evaluate each patch from each run on both the entire set of model misbehaviour and the sample it has seen during the patch generation, briefly checking the patch generalisability before RQ3 details it.

Additionally, we compare the performance of Arachne configured with our localisation method, against Arachne using RS and GL localisation, respectively, inspecting the effectiveness of different localisation methods for naturally emerging faults. Suppose our localisation method, BL, returns N_{avg} neural weights on average. With GL and RS, we generate an ordering as in RQ1, and select the top N_{avg} weights. To inspect whether Arachne can repair given misbehaviour of a DNN model, we report repair rate and break rate of the inputs (see details in Section 4.5).

3.2.3 RQ3. Repair Generalisability

Can Arachne generate a patch that generalises to unseen data?

The main purpose of Arachne is performing focused repair on a specific type of misbehaviour. This research question verifies that Arachne can repair target misbehaviour, and checks whether the generated patches are also effective for unseen inputs too. For this, we evaluate DNNs patched by Arachne using data unseen throughout the model training and the repair;³ We target the top 30 most frequent types of misclassifications and run Arachne on each misclassification type (i.e., inputs of class A being misclassified as class B) independently. To avoid any bias in neuron selection, we include all misclassified inputs of the given type in the set of negative inputs, I_{neg} . Subsequently, we evaluate DNN models before and after the patch against the same types of misclassification in unseen data.

3.2.4 RQ4. Balancing Behaviour

What is the impact of the balancing hyperparameter? The hyperparameter α in Arachne balances the amount of focus on fixing misbehaviour and on preserving correct behaviour. We evaluate the impact that different values of α have

³Section 4.3 details how we divided the test dataset for the repair and the evaluation.

on the repair outcome, using the most frequent type of misbehaviour for each dataset as an example. Using DNN classifiers studied in RQs 1 to 3, we investigate the impact of varying α .

3.2.5 RQ5. Comparison to State-of-the-Art

How does Arachne compare to the state-of-the-art in the context of targeted repair? We compare Arachne to the recent DNN repair technique, Apricot [17]. We fully⁴ train three DNN image classifiers using CNN architectures taken from Apricot and attempt to repair the three most frequent types of misbehaviour for each DNN classifier. In addition to these three CNNs from Apricot, we train and repair two additional image classifiers for two different datasets that were not studied with Apricot; here, we also target the top three most frequent types of misbehaviour of each classifier. While we reproduce most of the experimental settings of Apricot faithfully, we introduce a few changes for targeted repair. First, we modify how Apricot uses reduced Deep Learning Models (rDLMs). rDLMs are neural networks trained on a subset of the training data, to guide repair of a DNN based on their behaviour on individual inputs. We modify the rDLM-based weight adjustment so that it is initiated when a specific type of misbehaviour is discovered, instead of for every error. Second, we reduce the number of additional training epochs in Apricot from 20 to five, as the original setting resulted in inhibitive repair time. Third, we terminate Apricot when there are no improvements for 100 batches, again for time considerations. We repeat repairs using Arachne 30 times per model, while Apricot is executed once for each error type due to its long execution time.

In addition to Apricot, we compare against a *retraining baseline* that takes the fault localization results of the previous step, and uses newly provided images as training data to retrain the faulty weights. In this process, we freeze other non-targeted weights and focus the learning on the weights identified by our fault localization. Negative examples are weighted with α , equivalently to Arachne.

3.2.6 RQ6. Fairness Repair

We investigate whether Arachne can repair a more realistic and important DNN misbehaviour. Previous work has shown that commercial image classification APIs, trained with human faces, can reflect bias in the training data by showing accuracy gaps between genders [30]. We recreate this scenario by training a DNN model that classifies the gender of a given face image, using a pre-trained VGG architecture [37] and the Labelled Faces in the Wild (LFW) benchmark [26]. The classifier reflects the label imbalance in the dataset by showing lower class accuracy for the female gender. We apply Arachne to the trained model and see if it can successfully repair the label imbalance issue by patching the weights.

3.2.7 RQ7. Model Generalisation

To show that the use of Arachne is not restricted to CNN-based DNNs and further show that Arachne can operate on different domains than images, we evaluate Arachne on a completely different domain from previous RQs, namely an LSTM network trained on a text-based dataset. More specifically, we train an LSTM network model using the Twitter US Airline Sentiment dataset that contains tweets of travellers’ sentiments on airlines; the goal of this model is to predict the sentiments of these tweets. We then repair the most frequent misclassification of this model using Arachne, which, in this case, is predicting that the tweet is neutral instead of negative. One unique characteristic of the dataset is that it provides the confidence sentiment labelers had in their label, which may help show which inputs are being handled by Arachne. Hence, we first evaluate the repair performance of Arachne, then analyse the cases in which Arachne fails in conjunction with the confidence of manual labellers provided by the dataset we use.

4 Experimental Setup

This section describes the details of experimental setup.

4.1 Subjects

We use five well-known datasets to test the effectiveness of Arachne: *Fashion-MNIST*, *CIFAR-10*, *GTSRB*, *Labeled Faces in the Wild (LFW)*, and *Twitter US Airline Sentiment*.

⁴Note that subject DNNs were trained to meet a preset epoch number in Apricot [17], but our experiments found that the resulting models underfit the data under their setting.

4.1.1 Fashion-MNIST (FM)

FM [23] has been introduced to overcome the shortcomings of the widely studied image classification benchmark, MNIST [38]. Instead of hand-written digits, FM contains 60,000 training images and 10,000 test images of various fashion items. Each image is a grey-scale image of size (28,28), associated with one of ten labels. For this dataset, we train a neural network composed of one fully connected layer with 100 neurons, followed by a softmax layer with the cross-entropy loss function. We use this neural network as the DNN image classifier for the FM benchmark in RQs 1 to 4. Meanwhile in RQ5, we use the convnet network benchmark provided by the official FashionMNIST repository, which consists of two convolutional layers and two fully-connected layers, to evaluate Arachne under a more realistic scenario (in RQ5).

4.1.2 CIFAR-10 (C10)

CIFAR-10 contains 50,000 training images and 10,000 test images with ten different classes [24]; each image is an RGB image of size (32,32). For this dataset, we train a neural network with one convolutional layer with 16 filters, followed by a fully connected layer with 512 neurons and a softmax layer with cross-entropy as its loss function. In addition to this basic neural network, which we use to address RQs 1 to 4, we also employ three DNN architectures from Apricot, namely CNN1, 2, and 3 [17], in RQ5.

4.1.3 German Traffic Sign Recognition Benchmark (GTSRB)

The GTSRB benchmark contains 51,839 images (39,209 for training, 12,630 for test) of real-world German traffic signs labelled with 43 different classes; the size of each image varies from (15,15) to (250,250). For this dataset, we re-implement the CNN model proposed by *Team IDSIA* in the competition held at IJCNN 2011 [39]; we use this model to answer RQ5. Following how *Team IDSIA* handled the dataset, we resize each image to be (48,48) while using only the image within the bounding box that locates the traffic sign in the centre. For RQs 1 to 4, we use a relatively simple model: we train a neural network composed of one convolutional layer with 16 filters, followed by a fully-connected layer with 512 neurons and a softmax layer with a cross-entropy loss function.

4.1.4 Labelled Faces in the Wild (LFW)

The LFW dataset [26] contains 13,233 portrait images of various people, each labelled with the person’s identity (the original purpose of the LFW benchmark is to train face verification models). There are manually verified LFW gender labels made by Afifi & Abdelhamed [40], which we use to train a classifier that predicts the gender of an image. The labels show the dataset has about a 10:3 ratio between male and female faces: the bias towards male faces results in an accuracy gap between the two genders, which is similar to reported cases in commercial face recognition APIs [30]. Consequently, we believe the trained gender classification model exhibits a realistic fairness issue to evaluate Arachne on. We divide this dataset into 90% training set and 10% test set for experiments.

4.1.5 Twitter US Airline Sentiment

The Twitter US Airline Sentiment dataset is a collection of 14,640 real-world tweets regarding the six major US airlines; each input is in the *text* format, unlike previous datasets in which the inputs were images [27]. The data was gathered in February 2015, and each tweet was manually classified into three sentiment labels: positive, negative, or neutral. Along with sentiment labels, the human inspectors provided a ‘confidence score’, which is denoted as ‘airline sentiment confidence’, indicating how confident they were with their classification. The tweets in this dataset vary in length. Hence, we embed these tweets using GloVe word embedding vectors, pre-trained on 6 billion tokens, with 100 dimensions for each word vector [41]. We drop the name of airline during this preprocessing in order to avoid bias caused by mentioning a specific airline. We split this dataset, using 80% (11712) for the training and 20% (2928) for the repair and evaluation (10% for the holdout and 10% for the evaluation).

4.2 Algorithm Configuration

DE in Arachne requires the F and CR parameters to be set in advance, as well as the population size and the number of generations. We use the default parameters of the DE implementation in *scipy* [42]: CR is set to 0.7, and F is randomly sampled from (0.5, 1.0) for each generation. Following Piotrowski’s suggestion [34] that a population size of 100 suffices for problems with less than 30 dimensions, we set the population to 100, as the number of localised neural weights was mostly below 30. We set the maximum number of generations to 100, but terminate early if the best candidate solution does not change for ten consecutive generations.

4.3 Fitness Evaluation

Arachne requires both positive and negative input sets (i.e., I_{pos} and I_{neg}) to compute the fitness of candidate patches, as shown in Equation (3). In order to collect these two input sets, we divide the test set into halves, holding one half to collect the positive and negative inputs for simulating on-the-fly repair after the model deployment. We refer to this holdout input set as the *validation set*. We then evaluate the patches generated with these inputs using the other half of the test set, referred to as the *evaluation set*. While all RQs that involves the patch generation follow this dataset strategy, RQ2, whose aim is to inspect whether Arachne *can* repair given misclassified inputs, is an exception; here, we collect both positive and negative inputs from the entire test set instead of from only the validation set. For hyperparameter α , we set it to 10 as default, similar to previous work on Automatic Program Repair [43], but study the impact of varying its value in RQ4. For RQ6 and RQ7, we set it as 2 and 1, respectively; these values are obtained empirically.⁵

4.4 Model Training

We train DNN models until they do not improve after 100 consecutive epochs of training: among all the snapshots taken after each epoch, we choose the one with the highest test accuracy. This is done to prevent using under- or overfitted models. RQs 1 to 4 intend to validate the capability of Arachne in generating focused repair or ‘hot fixes’ of a certain type of misbehaviour, before evaluating its actual effectiveness. Hence, for these RQs, we train models with relatively low capacity, whose architectures are described in Section 4.1. These models obtain a training accuracy higher than 90%, even with their inherently limited capacity. More precisely, for FM, the training accuracy is 95.69%, and the test accuracy is 89.30%; for C10, the training accuracy is 90.74%, while the test accuracy is 72.49%. Lastly, for GTSRB, the training and test accuracy values are 99.99% and 97.41%.

We use Apricot [17] as the baseline for RQ5. Apricot uses the CIFAR-10 dataset. When comparing with Apricot, we use DNN architectures presented in their paper (CNN1, 2, 3) for a fair comparison. The training process is the same as the above. Each network achieves a training accuracy of 94.27%, 100%, and 99.56% and a test accuracy of 76.30%, 83.36%, and 76.38%, respectively. In addition to these three CNN models for CIFAR-10, we adopt two additional models, one for Fashion-MNIST and one for GTSRB, as explained in Section 4.1. These two models inherently have more capacity to learn compared to those used in RQs 1 to 4. We train these models in the same way as above: the train and test accuracy are 99.98% and 92.67% for FM, respectively, and for GTSRB, the accuracy is 100% for the training data and is 97.34% for the test data.

For RQ6, which is about fairness repair, we use a pre-trained VGG model as the starting point and tune the model for classifying face images in the LFW dataset. The training process is, again, the same as the above. The training accuracy of the tuned model is 99.67% overall and 99.90% and 98.89% for male and female images, respectively. For the test dataset, the overall accuracy is 96.15%; the accuracy for male images is 98.24%, and for the female images, it is 86.86%.

For RQ7, which addresses a different model that is not an image classifier, we train an LSTM-based sentiment classification model; this model consists of one LSTM layer with 256 units followed by a softmax layer with the cross-entropy as the loss. Following the same training strategy, we obtain the training accuracy of 91.50% and the test accuracy of 76.09% for this model.

4.5 Evaluation metric

We evaluate the localisation phase of Arachne by computing ROC-AUC of the localisation methods. Each localisation method returns an ordering of weights in terms of high suspiciousness. Subsequently, we plot the receiver operating characteristic (ROC) curve and calculate the area under the curve (AUC).

We evaluate the patch phase of Arachne from two perspectives: how successful patches are, and what their cost is. First, let us use the notation $c_{original} \rightarrow c_{predicted}$ to denote an input of class $c_{original}$ being misclassified as belonging to class $c_{predicted}$. For each error type $a \rightarrow b$, we evaluate the patch effectiveness using a set of misclassified inputs, $I_{neg} = I_{d,(a \rightarrow b)}$, in which d denotes the dataset these inputs are from (i.e., $d \in \{\text{validation (val)}, \text{evaluation (eval)}\}$). Consequently, $I_{d,(a \rightarrow b)} = \{i \in I_d \mid label_{GT}(i) = a \wedge label(i) = b\}$, where $label_{GT}(i)$ and $label(i)$ refer to the true and the predicted label of an input i .

⁵We suspect that smaller α worked better for RQ6 and RQ7 due to the number of unique labels in the dataset. Unlike FM, CIFAR-10, and GTSRB with either 10 or 43 labels, LFW (RQ6) and US Airline (RQ7) have only two and three labels and thereby have a higher risk of overfitting when using a larger α value.

To evaluate the degree of success Arachne repairs have, we define Repair Rate (RR) as shown in Equation (5). RR measures how many misclassifications Arachne corrected.

$$RR(I_{neg}) = \frac{[\# \text{ of patched input from } I_{neg}]}{|I_{neg}|} \quad (5)$$

In this equation, I_{neg} is a set of inputs initially misclassified by a model under repair; more specifically, I_{neg} is a set of misclassified inputs that Arachne aims to repair. We specify two specific repair rate metrics to measure both the degree of optimisation and patch generalisation: RR_{val} and RR_{eval} . $RR_{val} = RR(I_{neg})$ measures how well Arachne changed the classification results of images seen during repair. To measure generalisability, we need a few more definitions. Suppose we aim to reduce instances of $a \rightarrow b$ in a trained model. We apply Arachne to the model using $I_{val,(a \rightarrow b)}$ as the negative input set, I_{neg} . However, to evaluate how general the generated patch is, we calculate $RR_{eval} = RR(I_{eval,(a \rightarrow b)})$, which captures how many unseen misclassifications of type $a \rightarrow b$ Arachne has patched.

When we measure the cost of Arachne, we take two things into account: the side effects and efficiency. To measure the side effects of Arachne, we use two metrics: Break Rate (BR) defined in Equation (6), and accuracy per class. BR measures how many initially correct inputs were broken by the repair of Arachne, where I_{pos} is a set of inputs correctly classified by the original model.

$$BR = \frac{[\# \text{ of broken input from } I_{pos}]}{|I_{pos}|} \quad (6)$$

Accuracy/error per class breaks down the side effects by each class to provide a clearer view on how the model is affected. We use these two metrics flexibly between RQs, depending on their purposes. For efficiency, we simply measure the time it takes to repair the DNN.

4.6 Implementation & Environment

Arachne is implemented in Python version 3.6; our implementation of DE is extended from an example provided by DEAP [44]. DNN models, as well our baseline, are implemented using TensorFlow version 1.12.0 or PyTorch version 1.5.1. Code and model data are publicly available from <https://anonymous.4open.science/r/08bf3a0a-3d52-4f61-9c42-0aa2f9c94262/>. All time was measured on machines equipped with Intel Core i7 CPU, 32GB RAM, and NVidia TitanX GPU.

5 Results and Analysis

This section reports and discusses the results of the empirical evaluation, and answers the research questions described in Section 3.2.

5.1 RQ1: Localisation Effectiveness

Table 1 reports the ratio of inputs whose behaviour was changed by a mutated neural weight, namely, a fault, and the ROC-AUC values of different weight prioritisation schemes including ours. Figure 2 visualises the average ROC curves from 30 runs.

As shown in the Change Ratio column in Table 1, among the three data sets, Fashion-MNIST (FM) is the least affected by the perturbation of single neural weight, while GTSRB is most affected. Both the proposed localisation method, BL, and the baseline GL outperform RS in all three datasets, assigning the mutated neural weight a higher inspection priority. Compared to GL, BL obtains better performance in FM and C10, and performs similarly in GTSRB. The difference between the localisation performance is the most pronounced in FM, followed by C10 and GTSRB.

We suspect these variations in the performance difference are related to the impact of a single weight mutation on the model, which the Change Ratio reflects. As described in RQ1 in Section 3.2, we did not set an upper limit on the impact of a mutated neural weight on initial model behaviour: the perturbed neural weight can change more than 0.1% of the inputs. Indeed, the average change ratio varies from 0.61% to 4.94%, depending on the dataset. GL prioritises neural weights in descending order of their absolute gradient loss values. Hence, the more a single neural weight alone can affect the model behaviour, the greater its absolute gradient loss value is likely to be. As a result, compared to FM, GL obtains comparable performance to our method BL in C10 and GTSRB, the two datasets with a relatively higher change ratio. In GTSRB, both BL and GL obtain a close-to-perfect performance, obtaining 1.0 for rounded

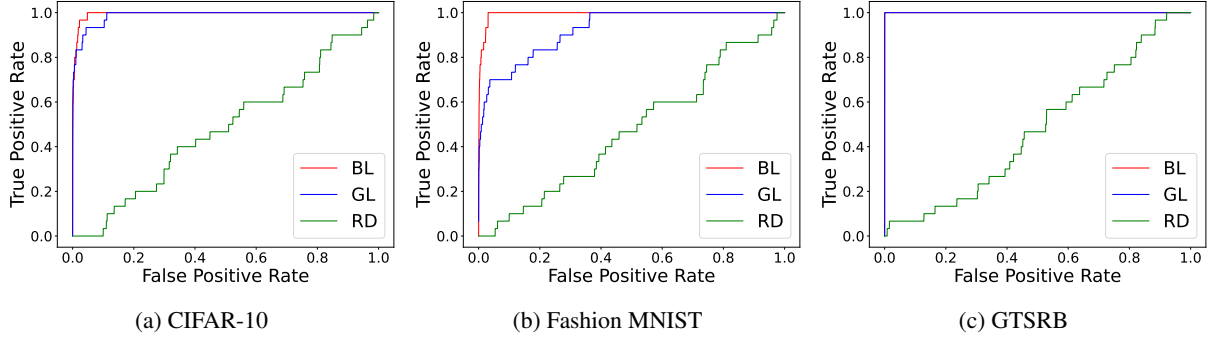


Figure 2: ROC curves for each localisation method. The lines indicate the median value over 30 runs; the shades represent the true positive value range including 90% of all runs.

ROC-AUC. We conjecture that this is due to the GTSRB dataset both containing more labels and fewer training data points. Together, these characteristics make the trained DNN model more sensitive toward small changes, especially when it obtains high accuracy. Indeed, the GTSRB model has a relatively higher change ratio and accuracy than the models of FM and C10.

We designed this localisation experiment to have a single faulty neural weight. However, in practice, multiple neural weights jointly play roles in the model misbehaviour, and therefore, the effect of a single neural will be less dramatic than in this experiment. Consequently, the notable performance of the proposed method, BL, in FM compared to other baselines implies that BL can effectively discriminate neural weights that cause the behavioural changes in the model.

Answer to RQ1: From the results, we conclude that the proposed BL method can effectively identify neural weights that are responsible for misbehaviour.

Table 1: The ratio of changed inputs (Change Ratio) and fault localisation accuracy (ROC-AUC) of the proposed method (BL), Gradient Loss (GL), and Random Selection (RS)

Subject	Change Ratio (%)			BL	GL	RS
	avg	min	max			
C10	0.61	0.12	2.43	0.9947	0.9880	0.4810
FM	0.17	0.1	0.42	0.9945	0.9237	0.4726
GTSRB	4.94	4.21	6.05	1.0000	1.0000	0.4774

5.2 RQ2: Patch Feasibility

Table 2 contains the average repair rate ($RR_{val} = RR(I_{neg})$) computed for the randomly selected ten percent of initial misclassified inputs, and the break rate (BR) from 30 repeated runs of Arachne with three different localisation methods; for the sake of simplicity, we denote RR_{val} as RR in Table 2. To summarise the overall repair cost, we further report the ratio of the post-patch model accuracy to the initial accuracy, R_{acc} , and the difference between the number of correct inputs before and after the patch, $\#_{diff}$, which is negative for performance degradation and positive for improvement, in the third column for each localisation method.

In order to mitigate potential bias from the input selection, at each run, we randomly select 10% of the total misclassified test inputs for Arachne to repair: 277, 107 and 125 inputs are selected for C10, FM, and GTSRB. Since the average number of localised neural weights is 11.6, 7.6, and 14.3 respectively for C10, FM, and GTSRB, GL and RS localisation methods were configured to select 12, 8, and 14 weights for C10, FM, and GTSRB.

With BL, Arachne successfully repairs on average 10.8% of the randomly selected misbehaviour in C10, while breaking the average 0.31% of the initially correct inputs (I_{pos}). RS barely modifies the network, and GL pays much more for the repair (i.e., breaks more inputs) than BL in C10. The R_{acc} column in Table 2 summarise the overall cost of the repair; compared to using BL for localisation, the misclassified inputs increase more than seven times in GL, as shown in the $\#_{diff}$ column (i.e., -14.4 for BL and -103.4 for GL).

For FM, Arachne with BL repairs 8.7% of the given misclassifications and breaks 0.7% of the initially correct ones, on average. Although Arachne repairs a little more when using GL (i.e., 9.5%), it also breaks more correct inputs. Similarly to what we observed in C10, RS repairs and breaks less than the other two. The accuracy decreases with all the three localisation methods but only slightly, retaining more than 99% of the initial accuracy.

For GTSRB, Arachne breaks less and repairs more or the same with BL compared to GL. For instance, Arachne repairs 13.0% of the given misclassifications using BL and 10.7%, a little less than BL, with GL. In terms of BR, BL outperforms GL by a small margin; the accuracy improves with BL and degrades with GL. Overall, the effect of Arachne is smaller in degree for GTSRB than for C10 and FM. We suspect this is due to the greater number of categories in the GTSRB dataset. As GTSRB has more categories (43) than the other two (10), the selected negative inputs of GTSRB likely include more various types of misclassifications than those of C10 and FM. Arachne repairs the given misbehaviour (i.e., misclassifications) of the model (I_{neg}) by directly modifying the neural weights specifically related to them. Hence, if there are multiple types of misbehaviour in I_{neg} , the overall repair effort of Arachne becomes distributed, resulting in less effective performance. Here, again, RS barely makes any meaningful change.

Arachne repairs the randomly chosen 10% of negative inputs instead of all of them in this experiment. Hence, we additionally compute the overall repair rate ($RR_{overall}$), RR calculated for the entire set of negative inputs instead of the randomly chosen 10%, to reduce the risk of drawing biased conclusions based on the selected inputs. Overall, $RR_{overall}$ shows a similar trend to that of $RR(=RR(I_{neg}))$, supporting the generalisability of our findings.⁶

While the underperformance of RS is what we expected (i.e., repair and break less), the results of GL require closer analysis. Since GL localises neural weights by following the gradient loss, we suspect that GL risks localising weights with small forward impact. Consequently, Arachne needs to make more dramatic changes to these weights in order to change the final classification outcome. In turn, this may have a greater effect on other inputs. BL does not suffer from these issues, and we believe that it can identify locations that have maximal patch effect with minimal modification. In summary, we conclude that Arachne can patch misclassifications, for which BL is the most effective localisation.

Answer to RQ2: Arachne can repair random misclassifications of a DNN while preserving most of the initially correct behaviour: repairing 10.8%, 8.7% and 13.0% of misclassifications in the test data, while breaking only 0.3%, 0.7 and 0.7% for C10, FM and GTSRB, respectively. The results suggest that BL is the most effective among the studied approaches in terms of balancing between repairing and preserving correct behaviour.

Table 2: The repair performance (RR and $RR_{overall}$) and cost (BR and R_{acc} (#diff)) of Arachne using bidirectional localisation method (BL), gradient loss (GL), and random selection (RS) over 30 runs.

Subject	BL			GL			RS		
	$RR / RR_{overall}$	BR	R_{acc} (#diff)	$RR / RR_{overall}$	BR	R_{acc} (#diff)	$RR / RR_{overall}$	BR	R_{acc} (#diff)
C10	0.108 / 0.074	0.031	0.998 (-14.4)	0.136 / 0.117	0.059	0.986 (-103.4)	0.017 / 0.007	0.001	1.001 (9.1)
FM	0.087 / 0.042	0.007	0.998 (-17.2)	0.095 / 0.047	0.008	0.998 (-20.0)	0.054 / 0.021	0.003	0.999 (-6.2)
GTSRB	0.130 / 0.069	0.007	1.001 (7.5)	0.107 / 0.068	0.008	0.999 (-8.1)	0.009 / 0.004	0.000	1.000 (4.8)

5.3 RQ3: Repair Generalisability

In RQ3, we aim to verify the generalisability of patches made by Arachne to fix specific misbehaviour, by reporting both RR_{val} and RR_{eval} . Table 3 reports the number of misclassifications repaired by Arachne in the validation and the evaluation dataset. The results suggest that a patch effective for a specific type of misclassifications in the validation data is also effective against the same type of misclassifications in the evaluation data. For example, Arachne successfully generates patches that fix 83% (86) of the 3 \rightarrow 5 misclassifications in the validation data of C10 on average; the same patches can repair 71% (74) of the same type of misclassified inputs in the evaluation data. Similarly, the patches that repair the average 75% (49) of 6 \rightarrow 0 misclassifications in FM can also repair the average 59% (38) of the same misclassifications in the evaluation data. In GTSRB, although the degree of repair is smaller in general, the patches that repair 53% (9) of the most frequent misclassifications (7 \rightarrow 0) can also fix 47% (8) of the same type of misclassifications in the evaluation data. The trend is similar in other misclassification types as well. Figure 3 visualises the overall results of Table 3.

Table 4 shows the overall generalisability results across all 30 studied misclassification types. We report the mean and standard deviation of the ratio between RR_{val} and RR_{eval} , aggregated over the top 1, 10, 20, and 30 most frequent misclassification types; for the top 1, we drop the standard deviation as there is only one ratio computed for the most frequent errors. The higher this ratio is, the more the patch generalises to unseen misclassifications in the test data. For

⁶RQ3 will detail the generalisability of Arachne.

this evaluation, we excluded the cases where Arachne cannot repair the given misclassifications at all ($RR_{val} = 0$); the numbers of such cases within the top 1, 10, 20 and 30 types are reported sequentially in Table 4, next to the project identifier.

Regarding the top 10 most frequent misclassification types in C10, Arachne keeps 88.48% of the repair rate (on average) with a standard deviation of only 0.1818 when moving from an observed dataset (validation) to an unobserved dataset (evaluation). Arachne shows similar results in FM, retaining 76.49% of the repair rate in the validation data for the top 10. In GTSRB, with both the average rate and the standard deviation of this ratio often greater than one, Arachne shows unstable performance compared to the other two. As can be inferred from Table 3, the number of misclassifications for each error type is generally smaller in GTSRB than in C10 and FM. This insufficient number of misclassification of GTSRB may explain the erratic performance we observed since it could increase the risk of generating close-to-random or over-fitted patches. The relatively low repair rate of Arachne in GTSRB in Table 3 further supports this suspicion. Nonetheless, Arachne still retains 66% of the repair rate for the most frequent type in GTSRB, promising its generalisability.

While this ratio tends to decrease as we gradually include less frequent misclassification types, for the top 30 most frequent types, Arachne successfully retains more than half of its repair rate against unseen data for all three datasets.

Answer to RQ3: The repair results for the 30 most frequent misclassification types show that Arachne can fix targeted misclassifications, and that the patches generated by Arachne generalise to unseen images.

Table 3: The number of inputs repaired by Arachne for the top ten most frequent misclassification types (Freq). The values in brackets are repair rates computed for each type.

Freq.		1	2	3	4	5	6	7	8	9	10
C10	fault	3 → 5	5 → 3	2 → 6	2 → 4	2 → 5	3 → 6	1 → 9	4 → 6	9 → 1	0 → 8
	val	86 (0.83)	35 (0.51)	28 (0.53)	34 (0.65)	30 (0.60)	27 (0.54)	35 (0.78)	22 (0.54)	31 (0.76)	28 (0.68)
	eval	74 (0.71)	39 (0.57)	25 (0.47)	36 (0.69)	24 (0.49)	25 (0.51)	22 (0.49)	20 (0.50)	32 (0.78)	23 (0.58)
FM	fault	6 → 0	6 → 2	2 → 4	0 → 6	4 → 2	6 → 4	2 → 6	4 → 6	9 → 7	6 → 3
	val	49 (0.75)	30 (0.64)	38 (0.88)	27 (0.64)	28 (0.74)	18 (0.51)	21 (0.68)	16 (0.67)	16 (0.80)	6 (0.38)
	eval	38 (0.59)	20 (0.43)	31 (0.72)	31 (0.74)	23 (0.62)	13 (0.38)	18 (0.58)	14 (0.58)	11 (0.55)	6 (0.38)
GTSRB	fault	7 → 8	41 → 9	30 → 31	23 → 31	17 → 38	6 → 5	18 → 21	19 → 31	6 → 42	2 → 3
	val	9 (0.53)	3 (0.19)	7 (0.50)	13 (1.00)	-1 (-0.08)	11 (0.92)	9 (0.82)	1 (0.09)	10 (0.91)	0 (0.00)
	eval	8 (0.47)	4 (0.25)	7 (0.54)	11 (0.92)	-2 (-0.15)	10 (0.91)	10 (0.91)	0 (0.00)	10 (1.00)	0 (0.00)

Table 4: The ratio between RR_{used} and RR_{eval} , computed for the top N frequent types of misclassification ($N \in \{1, 10, 20, 30\}$). The four numbers next to the project are the number of cases where Arachne completely fails to repair the given type of misclassifications for the top 1, 10, 20, and 30 frequent misclassification types.

top N	1	10		20		30	
		mean	std	mean	std	mean	std
C10 (0/0/0/0)	0.8571	0.8848	0.1818	0.7872	0.2850	0.7865	0.2935
FM (0/0/1/3)	0.7876	0.7649	0.2093	0.8062	0.2786	0.7058	0.3771
GTSRB (0/2/2/5)	0.6667	1.3735	1.6264	1.0188	1.1531	0.8282	1.0507

5.4 RQ4: Balancing Behaviour

Figure 4 shows changing class accuracies and overall accuracy when the hyperparameter α is configured differently. The results reveal that increasing α does result in more aggressive repairs, even if it comes at the cost of breaking more existing behaviour. For instance, when α is set to 10 in CIFAR-10 (making the impact of patches most notable), Figure 4 (b) shows that the patch increases the accuracy of the $C_{original}$ class by 40%p, at the cost of decreasing the accuracy of the $C_{predicted}$ class by 37%p. In contrast, when α is set to 1, Arachne neither fixes nor breaks as many inputs, repairing and breaking only 5.7% and 3.8% of the $C_{original}$ and $C_{predicted}$ classes. In Fashion MNIST, the overall trend is similar to that of CIFAR-10: with $\alpha = 10$, and for the evaluation data (in Figure 4 (d)), the accuracy of the $C_{original}$ increases by 18.81%p, whereas the accuracy of $C_{predicted}$ decreases by 24.4%. When $\alpha = 1$, the changes in both are within 1%p. The same goes for the GTSRB, despite on a far smaller scale. In the end, although the degree of the impact is different for each model and dataset, the trend of varying α is to be consistent across models (C10, FM, GTSRB) and datasets (validation and evaluation).

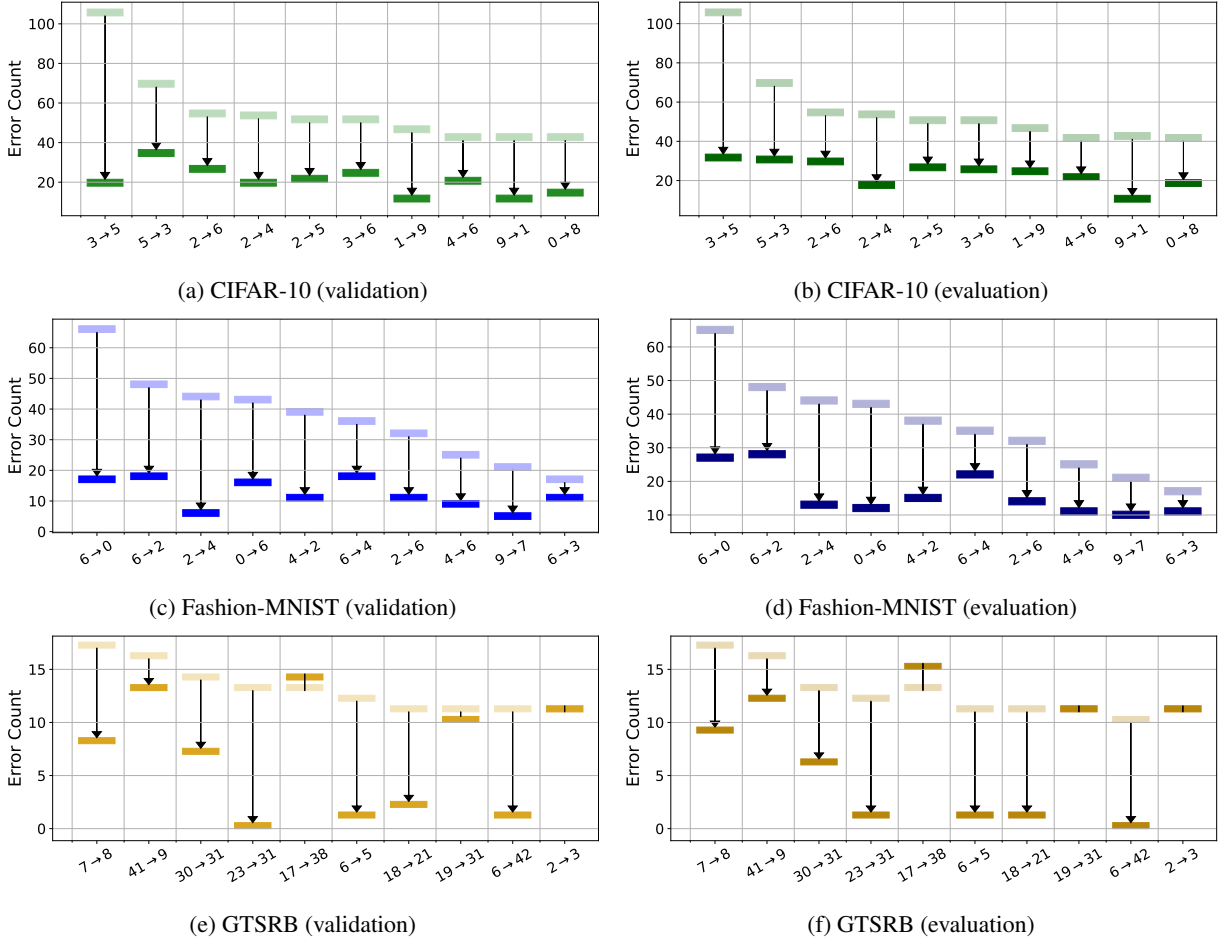


Figure 3: Visualisation of repair rate for the top 10 most frequent types of misclassification. The y-axis is the number of misclassifications that belong to each type. The light and dark bars denote the number of misclassifications before and after Arachne, respectively. The arrow from the light bar to the dark bar depicts the decrease in the number of misclassifications.

Answer to RQ4: We conclude that the hyperparameter α can effectively balance the behaviour of Arachne between aggressively pursuing repairs and conservatively retaining current behaviour.

5.5 RQ5: Comparison to the State-of-the-Art

Figure 5 compares the number of reduced misclassifications per type before and after repairs performed by Arachne and Apricot. Arachne successfully reduces the error count of targeted misclassification types across all models: across all fifteen combinations of five image classifier models and three misclassification types, the number of misclassifications is reduced by $RR_{eval} = 61.3\%$. Apricot, on the other hand, is less successful than Arachne, repairing only $RR_{eval} = 10.2\%$ of the targeted misclassifications on average.

Regarding the three image classifier models trained for CIFAR-10 (CNN1,2,3), Arachne successfully repairs the top three frequent misclassifications for all three models. In comparison, while Apricot does repair the misclassification for CNN3, the number of repaired errors is generally smaller than that of Arachne; for CNN1 and CNN2, Apricot does not report any change at all. We believe that the discrepancy in CNN1,2, and 3 between our results and those presented in Zhang and Chan [17] is due to the different initial settings: we use DNNs trained until their performance on the test set plateau, while Zhang and Chan [17] used DNNs trained until an arbitrary epoch threshold is met, making it possible that the DNNs in question were not fully trained. Meanwhile, in the FashionMNIST and GTSRB datasets that were not studied by Apricot, we observed that Arachne and Apricot perform similarly in generating targeted repairs

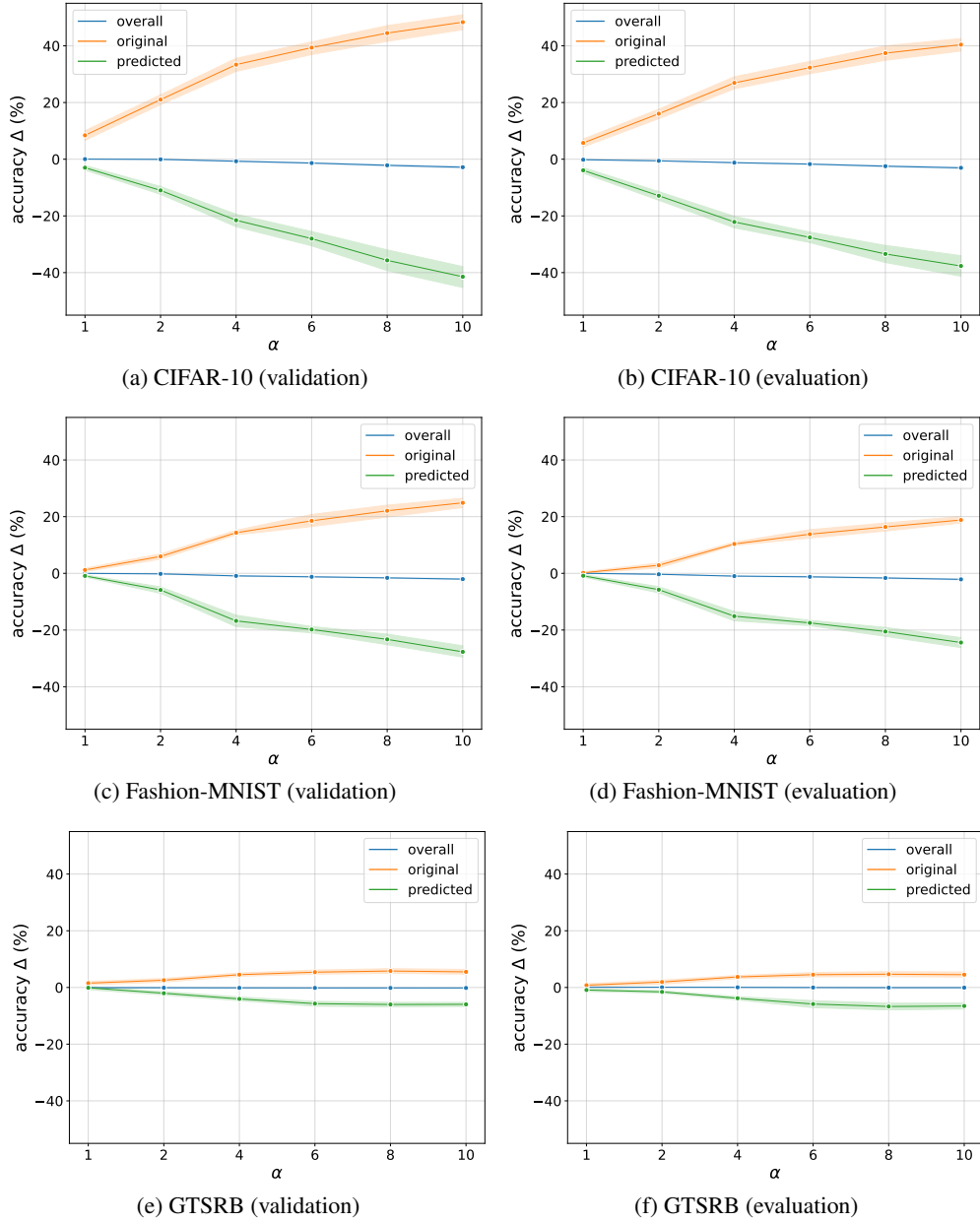


Figure 4: The Impact of Using Different α on the most frequent types of misclassifications. The x-axis is α used in Arachne, and the y-axis (Δ) is the corresponding degree of changes in accuracy. The orange and green circles signify the results of the original class and the predicted class of the target misclassified inputs, respectively. The blue circles show how the overall accuracy is affected.

in most cases except for the most frequent errors in FashionMNIST ($6 \rightarrow 0$), where Arachne outperformed Apricot; here, Arachne reduces the error count from 60 to near 20, whereas Apricot decreases it only to 50.

The high repair rate of Arachne may come at the cost of reduced overall accuracy, as Table 5 shows: when performing repair, Arachne may sacrifice overall accuracy to repair a specific type of misclassification. These results, combined with the higher repair rate, demonstrate that by construction Arachne is most useful as a targeted repair technique, and not as a retraining replacement for overall accuracy improvement. Meanwhile, we note that retraining barely changes performance; we believe this is due to the fact that while retraining uses gradient descent which can suffer when the optimization landscape is rough, Arachne uses the DE optimization technique which can overcome this issue.

Table 5: Average model accuracy of Arachne against the targeted version of Apricot and retraining (RT) for repairing frequent fault 1, 2, 3.

Model		Initial	Frequent Fault 1			Frequent Fault 2			Frequent Fault 3		
			Arac.	Apri.	RT	Arac.	Apri.	RT	Arac.	Apri.	RT
CNN1	Train	0.9427	0.9141	0.9427	0.9234	0.9066	0.9427	0.9234	0.9342	0.9427	0.9234
	Test	0.7630	0.7519	0.7630	0.7636	0.7531	0.7630	0.7636	0.7622	0.7630	0.7636
CNN2	Train	1.0000	0.9978	1.0000	1.0000	0.9956	1.0000	1.0000	1.0000	1.0000	1.0000
	Test	0.8336	0.8250	0.8336	0.8374	0.8223	0.8336	0.8374	0.8300	0.8336	0.8374
CNN3	Train	0.9956	0.9549	0.9996	0.9956	0.9692	0.9997	0.9852	0.9614	0.9995	0.9742
	Test	0.7638	0.7435	0.7871	0.7638	0.7525	0.7780	0.7568	0.7398	0.7759	0.7491
FM	Train	0.9948	0.9778	0.9997	0.9948	0.9841	0.9999	0.9948	0.9949	0.9999	0.9948
	Test	0.9267	0.9155	0.9250	0.9267	0.9212	0.9252	0.9267	0.9269	0.9242	0.9267
GTSRB	Train	1.0000	0.9999	1.0000	1.0000	0.9999	0.9989	1.0000	0.9999	1.0000	1.0000
	Test	0.9734	0.9731	0.9710	0.9734	0.9759	0.9678	0.9734	0.9732	0.9678	0.9734

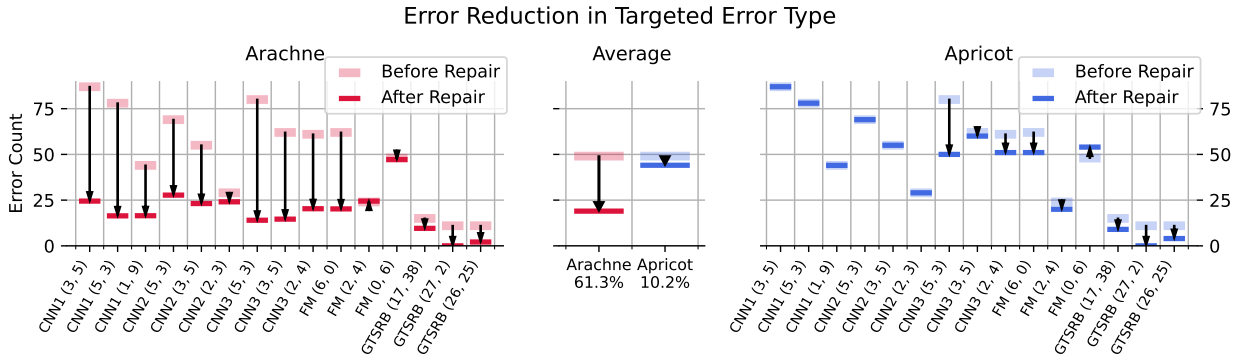


Figure 5: Number of wrong classifications of targeted type before and after repair, in Arachne and Apricot. The light bars represent the original errors of the targeted type, while the dark bars represent the average number of errors after the fix. These images were never seen by the neural network or the fixing algorithms. Retraining is omitted as it almost never changes performance. Lower is better.

We further compare the time it takes to perform repair on these DNNs: Table 6 shows the average repair time⁷ for the Frequent Fault 1. Arachne does not require any additional training epochs. On average, it terminates within ten minutes (i.e. two to seven mins.) for C10 and within three minutes for FM; For GTSRB, Arachne takes around average 23 minutes for the repair. In comparison, although we modified Apricot for early termination (reduced retraining epochs and early termination criterion: see Section 3.2.5), it takes nine (i.e., in GTSRB) to 86 times (i.e., in CNN3) longer than Arachne. These results further support the usefulness of Apricot in generating hot fixes that intend to quickly address errors of a certain type even at the cost of the others.

Table 6: Average patch time of Arachne against Apricot for repairing Frequent Fault 1.

model	Arachne	Apricot
CNN 1	2m 6s	1h 50m 12s
CNN 2	6m 8s	2h 16m 58s
CNN 3	6m 41s	9h 35m 52s
F-MNIST	2m 19s	1h 36m 4s
GTSRB	23m 19s	3h 29m 46s

Answer to RQ5: we answer RQ5 by noting that Arachne outperforms the state-of-the-art DNN repair technique, Apricot, both in terms of (targeted) repair rate and execution time.

⁷In this experiment, we measure the entire running time of Arachne, from the localisation to the repair

5.6 RQ6: Fairness Repair

RQ6 concerns a case study of a realistic fairness repair for DNNs. The gender classification DNN model we train using the LFW benchmark [26, 40] is biased due to the initial imbalance in the dataset: only 22.41% of the images in the LFW benchmark are those of female individuals. As a result, the error rate of the female class among the test data is eight times higher than that of the male class, reaching 13.2% against 2.7%, respectively.

We apply Arachne on the classifier, collecting I_{pos} and I_{neg} from the holdout test set (i.e., validation set): I_{pos} was all the images the classifier initially got correct, while I_{neg} consisted of 18 misclassified *female* images. Figure 6 shows how class accuracy changed when Arachne was applied. After repair, the female error per class for the test data changed to 11.5% and 5.6% for the female and male classes, respectively. The results suggest that we can use Arachne to address realistic fairness issues in DNN models, without resorting to data curation and retraining the model from scratch. We argue that issues such as model fairness can be an ideal application domain for Arachne, as quickly rebalancing the model may be necessary even if it degrades overall accuracy.

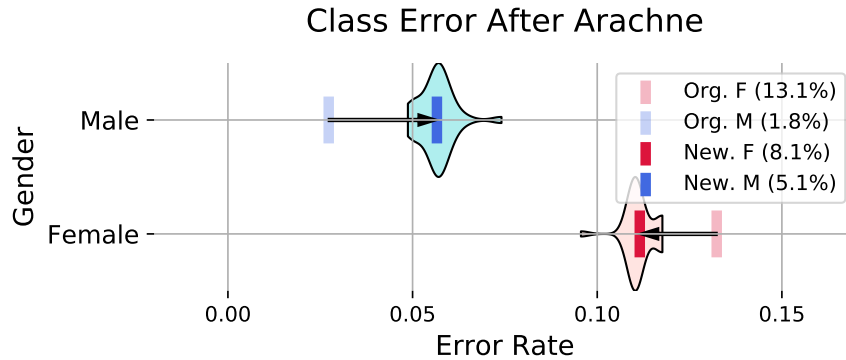


Figure 6: Class error before and after applying Arachne. Lower is better. The light contours represent the distribution of outcomes after 30 attempts to repair the neural network. The light bars indicate error by class before repair; the darker bars indicate average error after repair.

Answer to RQ6: Arachne can successfully repair DNNs with fairness issues, without additional data or retraining.

5.7 RQ7: Model Generalisation

RQ7 considers whether Arachne can generate fixes for different data modalities and different DNN architectures than those considered in RQ1 through RQ6, namely image-classifying CNNs; thus as explained previously, we evaluate Arachne on a text dataset (i.e., the tweet dataset) and with an LSTM DNN. The results of this analysis are presented in Figure 7a, in which we target the most frequent misclassification type, i.e., the misclassification of negative tweets as neutral. We posit that this type of misclassification poses a higher risk than other misclassifications, as the model will potentially miss out on important customer feedback. As the figure shows, Arachne reduces the prevalence of this error type not only on the validation set (which is the basis for the adjusting of the DNN weights), but also on the evaluation set which was unseen to Arachne, suggesting the changes that Arachne makes indeed generalize to new inputs. To reemphasize, these results demonstrate that Arachne can operate over diverse data modalities and DNN architectures.

We further compare the confidence of labeler in the assigned label (included in the dataset) between inputs whose classifications are corrected by the patch generated by Arachne, and inputs whose classifications remain incorrect even after the patch. The results of this comparison are provided in Figure 7b. We find that for the correctly patched inputs (i.e., inputs that were previously misclassified but now correctly classified), the confidence of human labelers is also high, whereas for the inputs Arachne fails to patch (i.e., inputs that continue to be misclassified despite the patch), the confidence of human labelers was low. These results suggest that Arachne repairs more worthwhile inputs, leaving only ambiguous cases misclassified. This provides further supporting evidence that Arachne can indeed generate worthwhile repairs for DNNs.

Answer to RQ7: Arachne can successfully repair completely different DNN architectures on different data modalities, suggesting the principles underpinning Arachne are general.

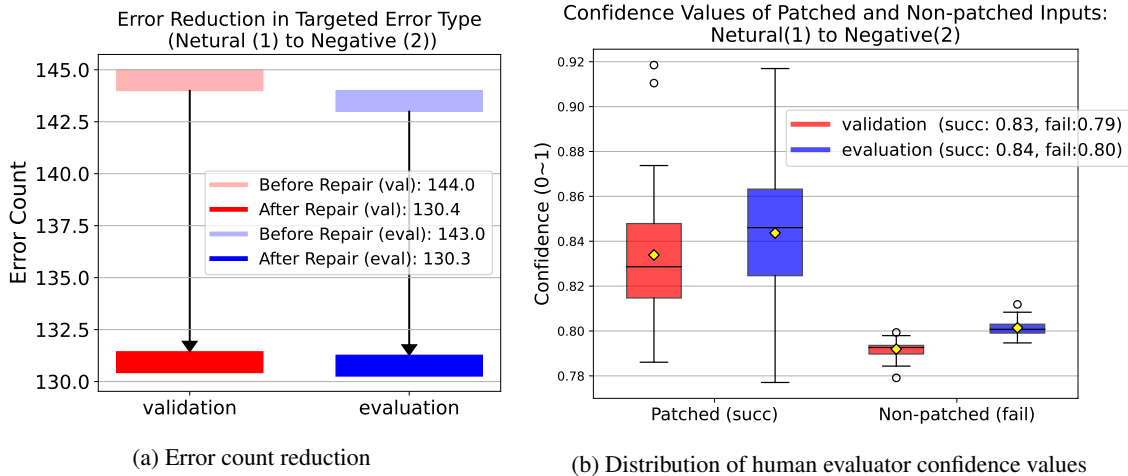


Figure 7: Results from repairing the LSTM classifier for Twitter US airline sentiment dataset. Blue and red denote the performance in the validation and evaluation set, respectively. The confidence values measure the confidence of human evaluator in their decision (positive/negative/neutral)

6 Threats to Validity

Threats to the internal validity of our empirical evaluation include the stochasticity of Arachne, correctness of data collection, and the correctness of our implementation of Apricot. Since both Arachne and the training of DNN models are stochastic processes, randomness may affect the observed results. We repeat experiments whenever possible and report results from multiple types of misclassification to avoid this. We use widely studied frameworks for evolutionary computation and deep learning to build our models as well as reconstruct Apricot: DEAP [44] and PyTorch, respectively. We also make our implementation available for public scrutiny.

Threats to external validity concern any factors that may limit the generalisation of our claim. CIFAR-10 and Fashion-MNIST are widely studied benchmarks for image classification. We also select GTSRB, the dataset that contains real-word traffic signs, to assess Arachne in a more practical case. We augment these with the fairness repair study based on the LFW benchmark, with results that support the generalisability of our claims. Lastly, we evaluate Arachne regarding the data modality and DNN architecture through the study other than the one on image classifier models, using a text dataset (i.e., the Twitter US airline sentiment dataset) and an LSTM DNN model. Additional studies can further strengthen the generalisability of our claims.

Threats to construct validity concern any potential misuse or misinterpretation of measured metrics. Our evaluation metrics are either standard classification metrics or simple count-based metrics with little room for misinterpretation.

7 Related Work

While there is a vast body of literature on techniques to improve the learning capability of DNNs [1], research on how to validate the quality of DNN models is relatively young. Existing literature focus on ways to reveal misbehaviour by using the combination of input synthesis and metamorphic oracles [10, 45]. Several test adequacy criteria have also been introduced and studied [10, 11, 12, 13].

Debugging or patching a DNN model has so far been formulated in the context of (re)training, i.e., continuing to learn until the correct behaviour is trained. Ma et al. proposed MODE [16], which uses Generative Adversarial Networks (GANs) to synthesise additional inputs that focus on the features relevant to the misbehaviour. These new inputs are used to retrain the DNN model under repair. Arachne, on the other hand, formulates the same problem as a direct Automated Program Repair (APR) for DNNs, and aims to solve it without generating or requiring additional data.

Zhang and Chan proposed Apricot [17], which trains multiple “reduced Deep Learning Models” (rDLMs) and incorporates them in an attempt to repair neural networks. Apricot trains multiple rDLMs, each of which have different classification results. By comparing the trained weights of rDLMs in these two different groups, Apricot readjusts the weights in the original DNN model towards correct rDLMs and applies additional training epochs. Apricot is the closest to Arachne among the existing techniques, in that it does not require additional training data. However, unlike

Apricot, Arachne does not rely on retraining at all. Our results suggest that the additional training may be the primary source of its time cost.

Another approach towards retraining focuses on which additional input to use for the retraining. DeepFault [46], for example, identifies *suspicious neurons* using a spectrum based approach, and perturbs existing input images in the direction that increases the activation value of the suspicious neurons. As a result, DeepFault can exploit the suspicious neurons (neurons that are correlated with DNN misbehaviour) to synthesize adversarial examples, which can subsequently be used for retraining. However, due to the way DeepFault synthesises new inputs based on domain constraints, its application is currently limited to Convolutional Neural Networks, whereas Arachne can be applied to any DNN architecture.

Arachne is heavily inspired by a class of APR techniques called Generate and Validate (G&V) [43, 18, 19, 47]. Concepts such as the use of positive and negative input sets, the use of fault localisation, and the use of metaheuristic search as the main driver of the repair are all inherited from existing G&V techniques. Regarding fault localisation, DeepFault [46] adopts the Spectrum Based Fault Localisation (SBFL) framework that has been successfully applied to coverage and test result information from traditional programs [31]. However, unlike Arachne that localises neural weights, DeepFault localises neurons. This coarser granularity of DeepFault requires an additional process to identify true-faulty neural weights among those connected to localised neurons for it to be used in the patch generation phase of Arachne. While treating all connected neural weights of faulty neurons as places to fix can be another option, this may lead to search space explosion. Further, DeepFault requires the user to set an arbitrary threshold for neural activation value to fit DNN activations into the program spectrum framework. Arachne does not require such hyperparameters, as its localisation depends on both gradient loss and forward impact of each neuron weight.

Regarding patching, like other G&V techniques, Arachne first *generates* a candidate patch and *validates* the patch by applying it and subsequently executing the patched model against the set of positive and negative inputs. However, due to the nature of DNNs, some fundamental differences exist. Both the program and the patch representation are numerical and continuous for Arachne. Unlike the fitness function of GenProg [43, 18] that only counts discrete step changes in the number of test cases that fail, Arachne can use the model loss as a continuous guidance towards repair. While this paper presents separate studies for accuracy and fairness repair, it is possible to repair a given DNN with respect to multiple objectives, for example by following the multi-objective formulation of APR introduced by ARJA [19]. This is possible because Arachne does not require any repair-specific modifications to the DNN.

8 Conclusion

We present Arachne, a search-based repair technique for DNN models. Arachne uses Differential Evolution to directly manipulate neural weight values, which are chosen by a novel localisation method. An empirical evaluation using widely studied four image benchmark datasets and one text dataset, suggests that Arachne can repair misbehaviour of DNN models while minimally disrupting existing correct behaviour. Furthermore, in targeted repair, Arachne outperforms a state-of-the-art method and tens of times faster. We also conduct a case study of realistic fairness repair: Arachne can repair a biased gender classification model by rebalancing the neural weights without requiring any additional data. The additional study with a text dataset and an LSTM-based model further supports the applicability of Arachne in different domains of datasets and DNN architectures. We believe that Arachne will facilitate the deployment of DNNs, acting as a fast adjustment technique capable of fixing both fairness and misclassification issues.

References

- [1] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, p. 436, 2015.
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Commun. ACM*, vol. 60, no. 6, pp. 84–90, May 2017. [Online]. Available: <http://doi.acm.org/10.1145/3065386>
- [3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [4] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, and B. Kingsbury, “Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 82–97, Nov 2012.
- [5] S. Jean, K. Cho, R. Memisevic, and Y. Bengio, “On using very large target vocabulary for neural machine translation,” in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long*

- Papers*). Beijing, China: Association for Computational Linguistics, Jul. 2015, pp. 1–10. [Online]. Available: <https://www.aclweb.org/anthology/P15-1001>
- [6] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *Proceedings of the 27th International Conference on Neural Information Processing Systems*, ser. NIPS’14. Cambridge, MA, USA: MIT Press, 2014, pp. 3104–3112.
 - [7] C. Chen, A. Seff, A. Kornhauser, and J. Xiao, “Deepdriving: Learning affordance for direct perception in autonomous driving,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2722–2730.
 - [8] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, “Multi-view 3d object detection network for autonomous driving,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1907–1915.
 - [9] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghahfoorian, J. A. van der Laak, B. van Ginneken, and C. I. Sánchez, “A survey on deep learning in medical image analysis,” *Medical Image Analysis*, vol. 42, pp. 60 – 88, 2017. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1361841517301135>
 - [10] K. Pei, Y. Cao, J. Yang, and S. Jana, “Deepxplore: Automated whitebox testing of deep learning systems,” in *Proceedings of the 26th Symposium on Operating Systems Principles*, ser. SOSP ’17. New York, NY, USA: ACM, 2017, pp. 1–18. [Online]. Available: <http://doi.acm.org/10.1145/3132747.3132785>
 - [11] Y. Tian, K. Pei, S. Jana, and B. Ray, “Deepest: Automated testing of deep-neural-network-driven autonomous cars,” in *Proceedings of the 40th International Conference on Software Engineering*, ser. ICSE ’18. New York, NY, USA: ACM, 2018, pp. 303–314. [Online]. Available: <http://doi.acm.org/10.1145/3180155.3180220>
 - [12] J. Kim, R. Feldt, and S. Yoo, “Guiding deep learning system testing using surprise adequacy,” in *Proceedings of the 41th International Conference on Software Engineering*, ser. ICSE 2019. IEEE Press, 2019, pp. 1039–1049.
 - [13] L. Ma, F. Juefei-Xu, F. Zhang, J. Sun, M. Xue, B. Li, C. Chen, T. Su, L. Li, Y. Liu, J. Zhao, and Y. Wang, “Deepgauge: Multi-granularity testing criteria for deep learning systems,” in *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, ser. ASE 2018. New York, NY, USA: ACM, 2018, pp. 120–131.
 - [14] M. Zhang, Y. Zhang, L. Zhang, C. Liu, and S. Khurshid, “Deeproad: Gan-based metamorphic testing and input validation framework for autonomous driving systems,” in *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, ser. ASE 2018. New York, NY, USA: ACM, 2018, pp. 132–142. [Online]. Available: <http://doi.acm.org/10.1145/3238147.3238187>
 - [15] S. Kang, R. Feldt, and S. Yoo, “Sinvad: Search-based image space navigation for dnn image classifier test input generation,” in *Proceedings of the International Workshop on Search Based Software Testing (SBST 2020)*, 2020.
 - [16] S. Ma, Y. Liu, W.-C. Lee, X. Zhang, and A. Grama, “Mode: Automated neural network model debugging via state differential analysis and input selection,” in *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2018. New York, NY, USA: ACM, 2018, pp. 175–186. [Online]. Available: <http://doi.acm.org/10.1145/3236024.3236082>
 - [17] H. Zhang and W. K. Chan, “Apricot: A weight-adaptation approach to fixing deep learning models,” in *Proceedings of the 34th IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE ’19. IEEE Press, 2019, pp. 376–387. [Online]. Available: <https://doi.org/10.1109/ASE.2019.00043>
 - [18] C. L. Goues, M. Dewey-Vogt, S. Forrest, and W. Weimer, “A systematic study of automated program repair: Fixing 55 out of 105 bugs for \$8 each,” in *Proceedings of the 34th International Conference on Software Engineering*, 2012, pp. 3–13.
 - [19] Y. Yuan and W. Banzhaf, “Arja: Automated repair of java programs via multi-objective genetic programming,” *IEEE Transactions on Software Engineering*, pp. 1–1, 2018.
 - [20] S. Saha, R. K. Saha, and M. R. Prasad, “Harnessing evolution for multi-hunk program repair,” in *Proceedings of the 41st International Conference on Software Engineering*, ser. ICSE ’19. Piscataway, NJ, USA: IEEE Press, 2019, pp. 13–24.
 - [21] R. Storn and K. Price, “Differential evolution – a simple and efficient heuristic for global optimization over continuous spaces,” *J. of Global Optimization*, vol. 11, no. 4, pp. 341–359, Dec. 1997. [Online]. Available: <https://doi.org/10.1023/A:1008202821328>
 - [22] S. Das and P. N. Suganthan, “Differential evolution: A survey of the state-of-the-art,” *IEEE Transactions on Evolutionary Computation*, vol. 15, no. 1, pp. 4–31, 2011.

- [23] H. Xiao, K. Rasul, and R. Vollgraf. (2017) Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms.
- [24] A. Krizhevsky, "Learning multiple layers of features from tiny images," University of Toronto, Tech. Rep., 2009.
- [25] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel, "The german traffic sign recognition benchmark: A multi-class classification competition," in *The 2011 International Joint Conference on Neural Networks*, 2011, pp. 1453–1460.
- [26] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," University of Massachusetts, Amherst, Tech. Rep. 07-49, October 2007.
- [27] Twitter us airline sentiment. [Online]. Available: <https://www.kaggle.com/datasets/crowdfLOWER/twitter-airline-sentiment>
- [28] S. Amershi, A. Begel, C. Bird, R. DeLine, H. Gall, E. Kamar, N. Nagappan, B. Nushi, and T. Zimmermann, "Software engineering for machine learning: A case study," in *Proceedings of the 41st International Conference on Software Engineering: Software Engineering in Practice*, ser. ICSE-SEIP 2019. IEEE Press, 2019, pp. 291–300.
- [29] Y. Roh, G. Heo, and S. E. Whang, "A survey on data collection for machine learning: a big data - AI integration perspective," *CoRR*, vol. abs/1811.03402, 2018. [Online]. Available: <http://arxiv.org/abs/1811.03402>
- [30] J. Buolamwini and T. Gebru, "Gender shades: Intersectional accuracy disparities in commercial gender classification," in *Conference on fairness, accountability and transparency*, 2018, pp. 77–91.
- [31] W. E. Wong, R. Gao, Y. Li, R. Abreu, and F. Wotawa, "A survey on software fault localization," *IEEE Transactions on Software Engineering*, vol. 42, no. 8, p. 707, August 2016.
- [32] W. Weimer, T. Nguyen, C. L. Goues, and S. Forrest, "Automatically finding patches using genetic programming," in *Proceedings of the 31st IEEE International Conference on Software Engineering (ICSE '09)*. Vancouver, Canada: IEEE, 16-24 May 2009, pp. 364–374.
- [33] J. Ronkkonen, S. Kukkonen, and K. V. Price, "Real-parameter optimization with differential evolution," in *2005 IEEE Congress on Evolutionary Computation*, vol. 1, 2005, pp. 506–513 Vol.1.
- [34] A. P. Piotrowski, "Review of differential evolution population size," *Swarm and Evolutionary Computation*, vol. 32, pp. 1 – 24, 2017. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S2210650216300268>
- [35] R. Storn and K. Price, "Differential evolution – a simple and efficient heuristic for global optimization over continuous spaces," *Journal of Global Optimization*, vol. 11, no. 4, pp. 341–359, 1997.
- [36] N. Humbatova, G. Jahangirova, G. Bavota, V. Riccio, A. Stocco, and P. Tonella, "Taxonomy of real faults in deep learning systems," in *The proceedings of the 42nd IEEE/ACM International Conference on Software Engineering*, ser. ICSE 2020, 2020, pp. 1110–1121.
- [37] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014. [Online]. Available: <http://arxiv.org/abs/1409.1556>
- [38] Y. LeCun, C. Cortes, and C. Burges, "Mnist handwritten digit database," *AT&T Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, vol. 2, 2010.
- [39] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel, "Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition," *Neural Networks*, no. 0, pp. –, 2012. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0893608012000457>
- [40] M. Afifi and A. Abdelhamed, "Afif4: deep gender classification based on adaboost-based fusion of isolated facial features and foggy faces," *Journal of Visual Communication and Image Representation*, vol. 62, pp. 77–86, 2019.
- [41] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532–1543. [Online]. Available: <http://www.aclweb.org/anthology/D14-1162>
- [42] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. Jarrod Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. Carey, Í. Polat, Y. Feng, E. W. Moore, J. Vand erPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and S. . . Contributors, "SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python," *Nature Methods*, vol. 17, pp. 261–272, 2020.

-
- [43] C. Le Goues, T. Nguyen, S. Forrest, and W. Weimer, “Genprog: A generic method for automatic software repair,” *IEEE Transactions on Software Engineering*, vol. 38, no. 1, pp. 54–72, Jan 2012.
 - [44] F.-A. Fortin, F.-M. De Rainville, M.-A. Gardner, M. Parizeau, and C. Gagné, “DEAP: Evolutionary algorithms made easy,” *Journal of Machine Learning Research*, vol. 13, pp. 2171–2175, July 2012.
 - [45] X. Huang, M. Kwiatkowska, S. Wang, and M. Wu, “Safety verification of deep neural networks,” in *Computer Aided Verification*, R. Majumdar and V. Kunčák, Eds. Cham: Springer International Publishing, 2017, pp. 3–29.
 - [46] H. F. Eniser, S. Gerasimou, and A. Sen, “Deepfault: Fault localization for deep neural networks,” in *Fundamental Approaches to Software Engineering*, R. Hähnle and W. van der Aalst, Eds. Cham: Springer International Publishing, 2019, pp. 171–191.
 - [47] M. Wen, J. Chen, R. Wu, D. Hao, and S.-C. Cheung, “Context-aware patch generation for better automated program repair,” in *Proceedings of the 40th International Conference on Software Engineering*, ser. ICSE ’18. New York, NY, USA: ACM, 2018, pp. 1–11. [Online]. Available: <http://doi.acm.org/10.1145/3180155.3180233>