# Evolution of Activation Functions for Deep Learning-Based Image Classification

Raz Lapid and Moshe Sipper*

sipper@bgu.ac.il
Department of Computer Science, Ben-Gurion University
Beer Sheva 84105, Israel

## ABSTRACT

Activation functions (AFs) play a pivotal role in the performance of neural networks. The Rectified Linear Unit (ReLU) is currently the most commonly used AF. Several replacements to ReLU have been suggested but improvements have proven inconsistent. Some AFs exhibit better performance for specific tasks, but it is hard to know a priori how to select the appropriate one(s). Studying both standard fully connected neural networks (FCNs) and convolutional neural networks (CNNs), we propose a novel, three-population, co-evolutionary algorithm to evolve AFs, and compare it to four other methods, both evolutionary and non-evolutionary. Tested on four datasets—MNIST, FashionMNIST, KMNIST, and USPS—coevolution proves to be a performant algorithm for finding good AFs and AF architectures.

## CCS CONCEPTS

• **Computing methodologies** → **Neural networks**; **Image processing**; • **Theory of computation** → **Evolutionary algorithms**.

## KEYWORDS

deep learning, activation functions, coevolution

## 1 INTRODUCTION

Artifical Neural Networks (ANNs), and, specifically, Deep Neural Networks (DNNs), have gained much traction in recent years and are now being effectively put to use in a variety of applications. Considerable work has been done to improve training and testing performance, including various initialization techniques, weight-tuning algorithms, different architectures, and more. However, one hyperparameter is usually left untouched: the activation function (AF). While recent work has seen the design of novel AFs [1, 22–24], the Rectified Linear Unit (ReLU) remains by far the most commonly

used one, mainly due to its overcoming the vanishing-gradient problem, thus affording faster learning and better performance.

AFs are crucial first and foremost because they transform a neural network from a simple linear algorithm into a non-linear one, thus allowing the network to learn complex mappings and functions. If the AFs of a neural network were removed, the whole architecture could be reduced to a linear operation on its input, with very limited use. Traditionally, the most commonly used AFs were sigmoid functions—such as logistic and hyperbolic tangent—which are bounded, differentiable, and monotonic. Such functions are prone to the vanishing-gradient problem [9]: when the function value is either too high or too low, the derivative becomes very small, and learning becomes very poor. To address this issue, different functions have been developed, the most prominent of which being the ReLU. Though it has proven highly performant, the ReLU is susceptible to a problem known as "dying" [11]: because all negative values are mapped to zero, there might be a scenario where a large number of ReLU neurons only output a 0 value. It might cause the entire network to "die", resulting in a constant function.

Different variants of ReLU, such as PReLU [32] and Leaky ReLU [31], have been proposed to address this issue. Although there has been much research towards creating different variants, ReLU remains the most commonly used AF. Highly popular architectures—including AlexNet, ZFNet, VGGNet, SegNet, GoogLeNet, SqueezeNet, ResNet, ResNeXt, MobileNet, and SeNet—have a common AF setup, with ReLUs in the hidden layers and a softmax function in the output layer [18].

This paper introduces a novel coevolutionary algorithm to evolve AFs for image-classification tasks. Our method is able to handle the simultaneous coevolution of three types of AFs: input-layer AFs, hidden-layer AFs, and output-layer AFs. We surmise that combining different AFs throughout the architecture may improve the network's performance. We compare our novel algorithm to four different methods: standard ReLU- or LeakyReLU-based networks, networks whose AFs are produced randomly, and two forms of single-population evolution, differing in whether an individual represents a single AF or three AFs. We chose ReLU and LeakyReLU as baseline AFs since we noticed that they are the most-used functions in the deep-learning domain.

In the next section we review the literature on evolving new AFs through evolutionary computation. Section 3 delineates the different methods used in this paper and our proposed coevolutionary algorithm. The experimental setup used to test the methods is delineated in Section 4. We present results and analyze them in Section 5. Finally, we offer concluding remarks in Section 6.

## 2 PREVIOUS WORK

There has been extensive research into using evolutionary algorithms for optimizing and improving ANNs and DNNs, which can be divided into two main categories: 1) topology optimization [13], and 2) hyperparameter and weight optimization [10]. Topology-optimization methods can be further divided into two groups: constructive and destructive. Constructive techniques start with a simple topology and gradually enhance its complexity until it meets an optimality requirement. Contrarily, destructive methods start with a complicated topology and gradually remove unneeded components.

One of the most well-known methods is NeuroEvolution of Augmenting Topologies (NEAT) [27]. This is an evolutionary algorithm-based method for developing neural networks. NEAT follows the constructive approach, with the evolutionary algorithm beginning with small, basic networks and gradually increasing their complexity over the generations. Finding intricate and complex neural networks is possible through this iterative approach. NEAT tries to modify network topologies and weights in an attempt to strike a balance between the fitness of developed solutions and their diversity.

[28] used a genetic algorithm to evolve neural-network architectures and weight initialization of deep CNNs for image-classification problems. Another work that used evolution to search for neural-network architectures is [30], in which they utilized a partial weight-sharing, one-shot neural architecture search framework that directly evolved complete neural-network architectures. [29] evolved neural-network architectures using a training-free genetic algorithm.

There are many more papers that aim to optimize a neural network's architecture without evolution, e.g., [8], wherein they introduced an automatic machine learning technique to optimise a CNN's architecture by predicting the performance of the network. Another interesting work is [7], in which they trained an over-parameterized network and then used a pruning criterion to prune the network.

Hyperparameter optimization is also a thriving field of research. The impact of hyperparameters on different deep-learning architectures has revealed complicated connections, with hyperparameters that increase performance in basic networks not having the same effect with more sophisticated topologies [6]. Selecting the right AF(s) for a specific task is crucial for the success of the neural network.

There are four commonly used methods for hyperparameter selection (including AFs): 1) manual search, 2) grid search, 3) random search, and 4) Bayesian optimization. Manual search simply refers to the user's manually selecting hyperparameters, a method often used in practice because it is quick and can result in sufficing results. Nonetheless, this technique makes it difficult to replicate results on new data, especially for non-experts. Grid search exhaustively generates sets of hyperparameters from values supplied by the user. This approach relies on the programmer's knowledge of the problem and produces repeatable results, but it is inefficient when exploring a large hyperparameter space. Grid search is commonly used because it is simple to set up, parallelize, and explore the whole search space (to some extent). Random search is similar

to grid search except that hyperparameter values are chosen at random from ranges specified by the user; it often works better than other methods [3]. Bayesian optimization is based upon using the information gained from a given trial experiment to decide how to adjust the hyperparameters for the next trial [26].

A recent work by [17] focused on evolving AFs for neural networks. Their work differs in several aspects from our novel coevolutionary algorithm:

(1) They used standard evolution, whereas we use coevolution.
(2) We deploy Cartesian genetic programming (CGP)—with its immanent bloat-control—rather than (bloat-susceptible) tree-based GP.
(3) We use evolved AFs in all network layers, as opposed to their work, which only used evolved AFs in the hidden layers.
(4) We use four benchmark image datasets in full, whereas they used a small subset of image datasets, and tabular datasets.
(5) We investigate two different architectures (fully connected neural network and convolutional neural network), with two different baseline AFs (ReLU and Leaky ReLU), while they used a randomly chosen architecture for each of the tested datasets.

A number of papers evolve AFs using genetic programming [2, 4, 5], but we have not found any using *coevolution* for this purpose. We believe coevolution is well-adapted to the problem at hand, given the different natures of AFs depending on their placement—input, hidden, or output layer.

We are interested in coevolving AFs for fixed-topology neural networks. The coevolutionary algorithm iteratively evolves better AFs, a technique that has received less attention so far within the scope of AFs.

## 3 METHODS

We compare five different methods for obtaining neural networks that perform a given task:

(1) Standard fully connected neural network (FCN) and standard convolutional neural network (CNN).
(2) Random FCN and random CNN.
(3) Single-population evolution, where an individual represents a single AF.
(4) Single-population evolution, where an individual represents three AFs: input-layer AF, hidden-layer AF, and output-layer AF.
(5) Coevolutionary algorithm, with three populations: a population of input-layer AFs, a population of hidden-layer AFs, and a population of output-layer AFs.

We used the following hyperparmaters for the learning algorithm in each of the methods: learning rate – 0.01, optimizer – Adam, batch size – $\frac{1}{10}$ the size of the training set. PyTorch's default initialization was used both for weights and biases, meaning that values were initialized from $U(-\sqrt{k}, \sqrt{k})$, where $k = \frac{1}{number\_of\_input\_features}$.

Below we delineate each of the above methods in more detail. We used the same number of candidate solutions in each of the different methods to conduct a fair comparison. Also, before moving on to describe the evolutionary-based algorithms, we provide a

**Table 1: FCN architecture.**

| Layer | Layer type | Number of neurons |
|---|---|---|
| 1 (input) | Fully Connected | 784 |
| 2-6 (hidden) | Fully Connected ($\times$ 5) | 32 |
| 7 (output) | Fully Connected | 10 |

**Table 2: CNN architecture.**

| Layer | Layer type | No. channels | Filter size | Stride |
|---|---|---|---|---|
| 1 | Convolution | 32 | $3 \times 3$ | 1 |
| 2 | Max Pooling | N/A | $3 \times 3$ | 2 |
| 3 | Convolution | 64 | $3 \times 3$ | 1 |
| 4 | Max Pooling | N/A | $3 \times 3$ | 2 |
| 5 | Dropout ($p = 0.5$) | N/A | N/A | N/A |
| 6 | Fully Connected | 1600 | N/A | N/A |
| 7 | Fully Connected | 1000 | N/A | N/A |
| 8 | Fully Connected | 10 | N/A | N/A |

brief description of Cartesian genetic programming, the flavor of evolutionary algorithm we used herein.

## 3.1 Standard neural networks

A standard FCN is typically composed of several building blocks that are chained sequentially, starting with data input and ending with an output (upon which a loss function is computed). In the *Standard-FCN* architecture we used, the non-linear functions torch.nn.ReLU or torch.nn.LeakyReLU follow each fully connected layer (and before the final softmax layer). In the *Standard-CNN* architecture we used, torch.nn.ReLU or torch.nn.LeakyReLU follow each convolution layer and the last fully connected layer (before the final softmax layer). ReLU and Leaky ReLU thus form our baseline AFs for comparison.

In our setting (for all methods), a network requires three AFs (which can be identical—or not): an input-layer AF, a hidden-layer AF, and an output-layer AF. For FCNs, the first AF is applied after the first torch.nn.Linear layer, the second AF is applied after each of the hidden torch.nn.Linear layers, and the third AF is applied after the last torch.nn.Linear layer (before the softmax layer; see Table 1). For CNNs, the first AF is applied after the first torch.nn.Conv2D layer, the second AF is applied after the second torch.nn.Conv2D layer, and the third AF is applied after the first torch.nn.Linear layer (see Table 2).

## 3.2 Random neural networks

Random search is a class of numerical optimization methods that do not consider the problem's gradient, allowing them to be utilized on functions that are neither continuous nor differentiable. Herein, we randomly created a set of candidate (AF) solutions, with each one comprising three ordered AFs. We generated *max_iter*\*3 candidate solutions for fair comparison with CGP (*max_iter*—see Table 3), by repeatedly creating a random initial generation through CGP. Each candidate solution was then trained over the $\text{train}_1$ set for 3 epochs[1] and its score was computed as the accuracy over the $\text{train}_2$

---
[1]3 epochs proved a good tradeoff between evolution-driving fitness and runtime.

**Table 3: Hyperparameters used by CGP.**

| Parameter | Description | Value |
|---|---|---|
| *operators* | list of primitives | see Table 4 |
| *n_const* | number of symbolic constants | 0 |
| *n_rows* | number of rows in the code block | 5 |
| *n_columns* | number of columns in the code block | 5 |
| *n_back* | number of rows to look back for connections | 5 |
| *n_mutations* | number of mutations per offspring | 3 |
| *mutation_method* | specific mutation method | point |
| *max_iter* | maximum number of generations | 50 |
| *lambda* | number of offspring per generation | 4 |
| *f_tol* | absolute error acceptable for convergence | 0.01 |
| *n_jobs* | number of jobs | 1 |

set. Then, we pick the best (top-accuracy) candidate solution for comparison.
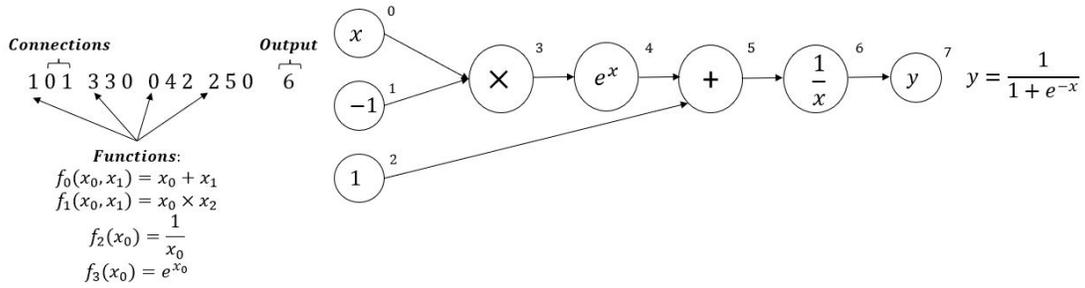
## 3.3 Cartesian genetic programming

Cartesian genetic programming (CGP) is an evolutionary algorithm wherein an evolving individual is represented as a two-dimensional grid of computational nodes—often an a-cyclic graph—which together express a program [14]. It originally grew from a mechanism for developing digital circuits [16]. An individual is represented by a linear genome, composed of integer genes, each encoding a single node in the graph, which represents a specific function. A node consists of a function, from a given table of functions, and connections, specifying where the data for the node comes from.

A sample individual is shown in Figure 1. There are three hyperparameters that define the connectivity and dimensionality of the encoded architecture: 1) number of rows, 2) number of columns, and 3) number of rows to look back for connections—this limits the columns from which a node can acquire its inputs. This representation is simple, flexible, and convenient for many problems. Typically, evolution begins with a population of randomly selected candidate solutions. In our case, because we are evolving AFs, an individual in a population represents an AF. Driven by a fitness function and using stochastic modification operators, CGP is able to produce successively better models in an iterative manner.

The CGP library used in this paper [21] evolves the population in a $(1 + \lambda)$-manner, i.e., in each generation it creates $\lambda$ offspring (we used the default $\lambda = 4$) and compares their fitness to the parent individual. The fittest individual carries over to the next generation; in case of a draw, the offspring is preferred over the parent. Tournament selection is used, with tournament size $k = |population|$, single-point mutation, and no crossover. For more details, see [15, 21]. Table 3 delineates the hyperparameters used by the CGP algorithm, and Table 4 shows the primitive set.

## 3.4 Evo-Single: Single-population evolution

Evo-Single is a single-population evolutionary algorithm, wherein an individual represents a single AF used throughout the whole network (as explained in Section 3.1). We start with the ReLU AF or the Leaky ReLU AF. Then we iterate, mutating the AF at every iteration of the CGP algorithm. Each individual is then trained over the $\text{train}_1$ set for 3 epochs and its fitness is computed as the accuracy over the $\text{train}_2$ set. After $max\_iter * 3 = 150$ iterations, we select the individual whose performance on the $\text{train}_2$ set was

$$y = \frac{1}{1+e^{-x}}$$

**Functions**:
$$f_0(x_0, x_1) = x_0 + x_1$$
$$f_1(x_0, x_1) = x_0 \times x_2$$
$$f_2(x_0) = \frac{1}{x_0}$$
$$f_3(x_0) = e^{x_0}$$

**Figure 1: A sample CGP individual, with 3 inputs—$x, -1, 1$—and 1 output—$y$. The genome consists of 5 3-valued genes, per 4 functional units, plus the output specification (no genes for the inputs). The first value of each 3-valued gene is the function's index in the lookup table of functions (bottom-left), and the remaining two values are parameter nodes. The last gene determines the outputs to return. In the example above, with $n_i$ representing node $i$: node 3, gene 101, $f_1(n_0, n_1) = n_0 \times n_1$; node 4, gene 330, $f_3(n_3) = e^{n_3}$ (unary function, 3rd gene value ignored); node 5, gene 042, $f_0(n_4, n_2) = n_4 + n_2$; node 6, gene 250, $f_2(n_5) = \frac{1}{n_5}$; node 7, output node, $n_6$ is the designated output value. The topology is fixed throughout evolution, while the genome evolves.**
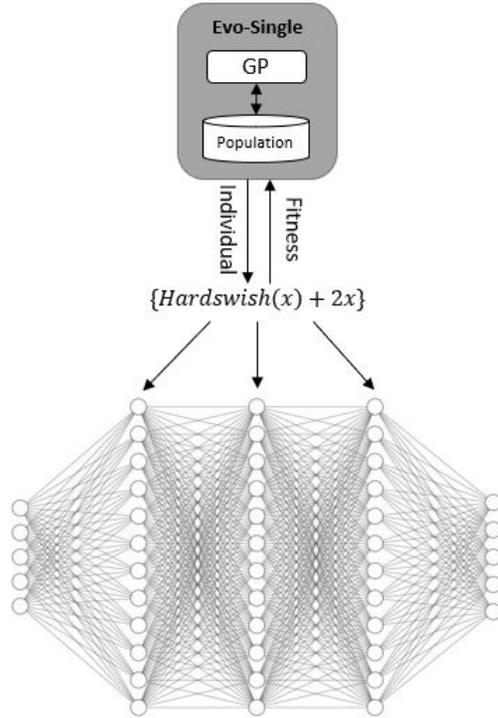
**Table 4: Set of primitives used by CGP. Left: Standard AFs. Right: Mathematical functions.**

| AF | Expression | AF | Expression |
|---|---|---|---|
| ReLU | $f(x) = max(x, 0)$ | Max | $f(x, y) = max(x, y)$ |
| Tanh | $f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ | Min | $f(x, y) = min(x, y)$ |
| Leaky ReLU | $f(x < 0) = negative\_slope \cdot x$ | Add | $f(x, y) = x + y$ |
| | $f(x \geq 0) = x$ | | |
| ELU | $f(x \leq 0) = (exp(x) - 1)$ | Sub | $f(x, y) = x - y$ |
| | $f(x > 0) = x$ | | |
| HardShrink | $f(x < -\lambda) = x$ | Mul | $f(x, y) = x \cdot y$ |
| | $f(x > \lambda) = x$ | | |
| | $f(-\lambda \leq x \leq \lambda) = 0$ | | |
| CELU | $f(x) = max(0, x) + min(0, \alpha(e^{\frac{x}{\alpha}} - 1)$ | | |
| Hardtanh | $f(x > 1) = 1$ | | |
| | $f(x < -1) = -1$ | | |
| | $f(-1 \leq x \leq 1) = x$ | | |
| Hardswish | $f(x \leq -3) = 0$ | | |
| | $f(x \geq 3) = -1$ | | |
| | $f(-3 < x \leq 3) = x \cdot \frac{(x+3)}{6}$ | | |
| Softshrink | $f(x > \lambda) = x - \lambda$ | | |
| | $f(x < -\lambda) = x + \lambda$ | | |
| | $f(-\lambda \leq x \leq \lambda) = 0$ | | |
| RReLU | $f(x \geq 0) = x$ | | |
| | $f(x < 0) = \alpha \cdot x$ | | |
| | ($\alpha$ is uniformly, randomly sampled) | | |

best. This individual is then compared against the other methods. The fitness-computation scheme is shown in Figure 2.

### 3.5 Evo-Triple: Single-population evolution

Evo-Triple is a single-population evolutionary algorithm, wherein an individual comprises an array of 3 AFs, the first being the input-layer AF, the second being the hidden-layer AF, and the third being the output-layer AF. The AFs are used throughout the network as explained in Section 3.1. We start with 3 ReLU AFs or 3 Leaky ReLU AFs. Then we iterate, mutating each of the AFs at every iteration of the CGP algorithm. Each individual is then trained over the $train_1$ set for 3 epochs and its fitness is computed as the accuracy over the $train_2$ set. After $max\_iter * 3 = 150$ iterations, we select the individual whose performance on the $train_2$ set was best. This individual is than compared against the other methods. As with Evo-Single, training uses the $train_1$ set, and fitness is computed



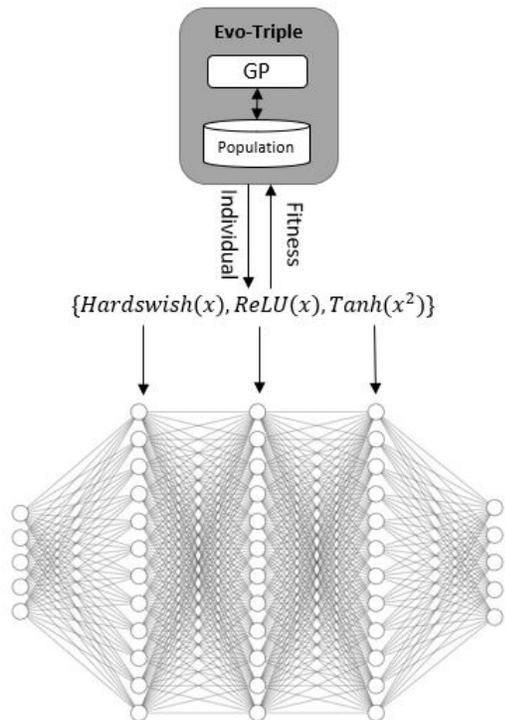**Figure 2: Evo-Single: Fitness computation of an individual in the population. The same AF is used throughout the network, in this example $Hardswish(x) + 2x$.**

over the $train_2$ set. The fitness-computation scheme is shown in Figure 3.

### 3.6 Coevo: Three-population coevolution

Coevolution refers to the simultaneous evolution of two or more species with coupled fitness [20]. Coevolution can be mutualistic,

**Figure 3: Evo-Triple: Fitness computation of an individual in the population, which comprises 3 AFs, in this example: $\{Hardswish(x), ReLU(x), Tanh(x^2)\}$.**

**Table 5: Datasets. Note: $1 \times h \times w$ represents a single-channel (greyscale) image of height $h$ and width $w$.**

| Dataset | Images | Classes | Training | Test |
|---|---|---|---|---|
| MNIST | $1 \times 28 \times 28$ | 10 | 60,000 | 10,000 |
| KMNIST | $1 \times 28 \times 28$ | 10 | 60,000 | 10,000 |
| FashionMNIST | $1 \times 28 \times 28$ | 10 | 60,000 | 10,000 |
| USPS | $1 \times 16 \times 16$ | 10 | 7291 | 2007 |

There might arise a concern over the differentiability of evolved AFs (PyTorch performs automatic differentiation—indeed this feature is one of its main strengths). However, in practice, this does not pose a problem since the number of points at which most functions are non-differentiable is usually (infinitely) small. For example, ReLU is non-differentiable at zero, however, in practice, it is rare to stumble often onto this value, and when we do, an arbitrary value can be assigned (e.g., zero), with little to no effect on performance. Our experiments have not shown differentiability to be of any concern.

## 4 EXPERIMENTAL SETUP

We experimented with four supervised classification datasets—MNIST, FashionMNIST, KMNIST, USPS—using PyTorch's [19] default data loaders, with pixel values normalized between $[0, 1]$ in the preprocessing step. The architectures used in the experiments are fixed and delineated in Tables 1 and 2. Full information about the datasets is shown in Table 5.

We used the Torchvision library—an accessible, open-source, machine-vision package for PyTorch, written in C++ [12]. We used our GPU cluster[2], performing 20 replicate runs for each experimental setup. We split the data three ways into $train_1$ (60%, used to "locally" train the networks), $train_2$ (25%, used to "globally" evolve the networks), and test (15%) sets, using the training sets for the different automated approaches, while saving the test set for post-evaluation.

As noted above, we compared 5 different approaches: Standard-FCN/CNN, Random-FCN/CNN, Evo-Single, Evo-Triple, and Coevo. We ran multiple replicates of four kinds, each involving the execution and comparison of five different methods:

(1) FCN, ReLU: Run Standard-FCN with ReLU, Random-FCN, Evo-Single, Evo-Triple, Coevo.
(2) FCN, Leaky ReLU: Run Standard-FCN with Leaky ReLU, Random-FCN, Evo-Single, Evo-Triple, Coevo.
(3) CNN, ReLU: Run Standard-CNN with ReLU, Random-CNN, Evo-Single, Evo-Triple, Coevo.
(4) CNN, Leaky ReLU: Run Standard-CNN with Leaky ReLU, Random-CNN, Evo-Single, Evo-Triple, Coevo.

The score of a network at any given phase was its classification accuracy over said phase's dataset. Scoring a standard network (either FCN or CNN) was done as follows: train the network using as training set both the $train_1$ and $train_2$ sets; the trained network's final score for comparison purposes was its performance over the test set. Scoring the other types of networks (random, evolved, and coevolved) was done as follows: train the network using the $train_1$

parasitic, or commensalistic [25]. In mutualistic coevolution, two or more species reciprocally affect each other beneficially. Parasitic coevolution is a competitive relationship between species. Commensalistic coevolution is a relationship between species where one of the species gains benefits while the other neither benefits nor is harmed. In this paper we focus on mutualistic, or cooperative coevolution, which is based on the collaboration of several independently evolving species to solve a problem. An individual's fitness is determined by its ability to cooperate with members of the other species.

Our coevolutionary algorithm comprises three separate populations: 1) input-layer AFs, 2) hidden-layer AFs, 3) output-layer AFs. Combining three individuals—one from each population—results in an AF architecture that can be evaluated. For each evolved individual in each population—an AF—we evaluated its fitness by choosing the two best (top-fitness) individuals of the two other populations from the previous generation. Each population evolved for $max\_iter = 50$ iterations (generations). We then created a neural network composed of the three AFs in the respective input, hidden, and output layers. The scheme for fitness computation is depicted in Figure 4. As with Evo-Single and Evo-Triple, training uses the $train_1$ set, and fitness is computed over the $train_2$ set. The fitness-computation scheme is shown in Figure 4. The coevolutionary algorithm is described in Algorithm 1.

---

[2]Nvidia GPU cards: 19 RTX 3090, 56 RTX 2080, 52 GTX 1080, 2 Titan XP, 4 P100.
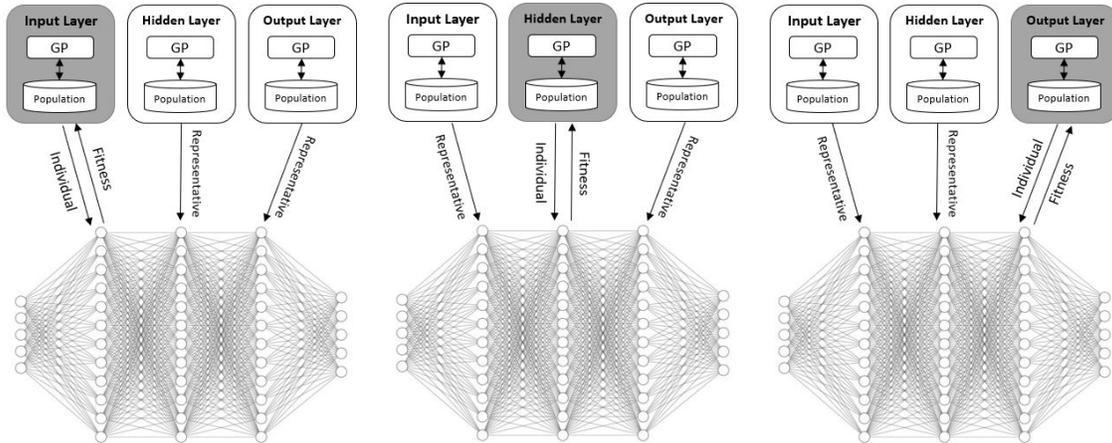
**Figure 4: Fitness computation of an individual in Coevo requires cooperation from the other two populations.**

---

**Algorithm 1** Coevolutionary Algorithm: Main Components.

1: **function** CREATE_NETWORK(*individual*, *best_input*, *best_hidden*, *best_output*, *pop*)  # *pop* $\in \{input, hidden, output\}$
2:     **if** $pop == input$ **then**
3:         *model* ← NEURALNETWORK(*individual, best_hidden, best_output*)
4:     **else if** $pop == hidden$ **then**
5:         *model* ← NEURALNETWORK(*best_input, individual, best_output*)
6:     **else** # $pop == output$
7:         *model* ← NEURALNETWORK(*best_input, best_hidden, individual*)
    **return** model

8: **function** NEURALNETWORK(*input_af*, *hidden_af*, *output_af*)
9:     *model* ← Create a neural network with the given AFs, and input/output dimensions determined by dataset
10:     **return** *model*

11: **function** FITNESS(*model*, *train₁*, *train₂*)
12:     train model for 3 epochs with cross entropy loss and mini-batches of size $|train_1|/10$
13:     **return** *accuracy score* over *train₂*

---

set; use as fitness the score over the $train_2$ set; the network's final score for comparison purposes was its performance over the test set. Note that the random and evolutionary algorithms never saw the test set—only the $train_1$ and $train_2$ sets were provided; solely the best individual returned was ascribed a test score based on the test set. In the next section, wherein we describe our results, values shown are test scores (i.e., over test sets).

Algorithm 2 delineates the experimental setup in pseudo-code format.

Figure 5 depicts an overview of the coevolutionary experimental setup. The code is available at https://github.com/razla.

## 5 RESULTS AND ANALYSIS

Our experimental results are shown in Table 6. The mean of the highest per-replicate test accuracy score over all replicates, for each method, is reported as the best score for that method.

Of the 16 groups of runs, Coevo won 7, Evo-Triple – 6, Random – 2, Standard – 1, and Evo-Single – 0 (wins are sometimes by a small margin but this is usual in this field). The standard FCN or

---

**Algorithm 2** Experimental Setup

**Input:**
    *dataset* ← dataset to be used
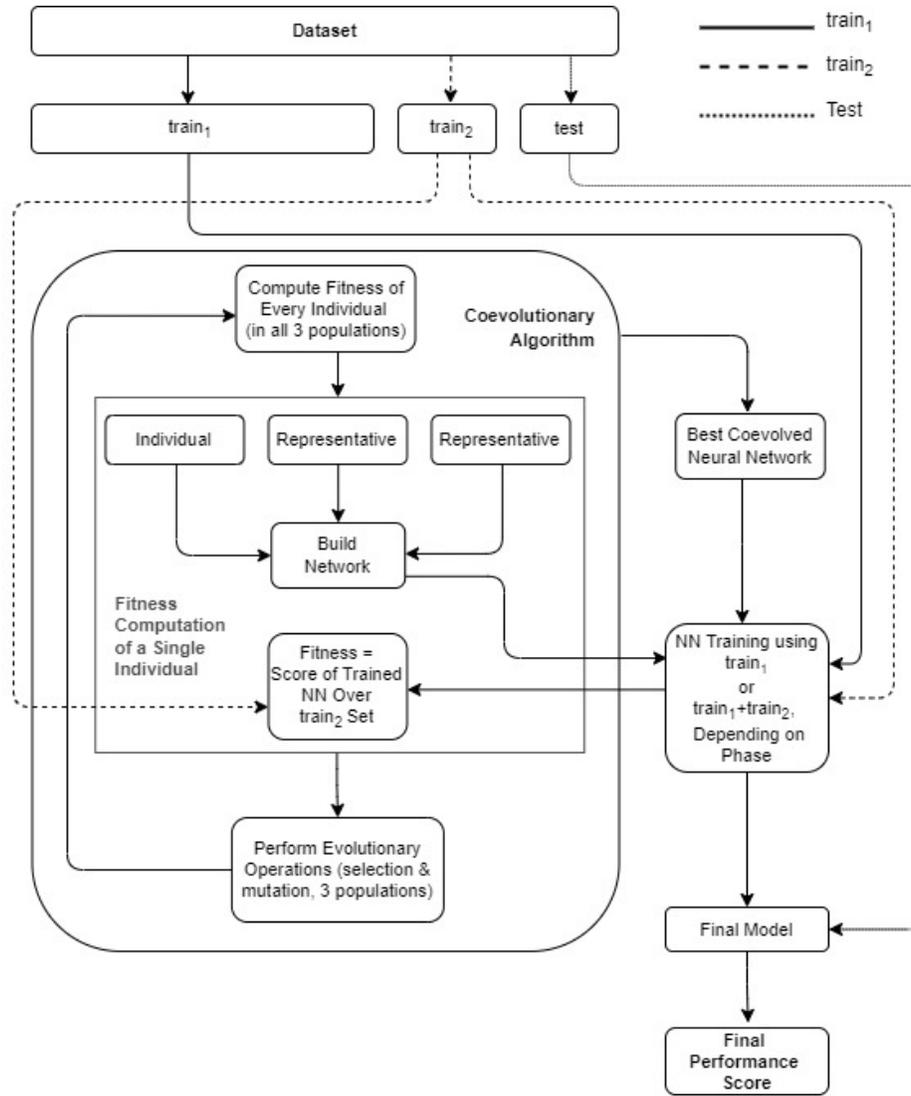    *replicates* ← number of replicates to run
**Output:**
    Performance measures (over test sets)

1: **for** *rep* in {1…*replicates*} **do**
2:     Shuffle and split the dataset into *train₁* set, *train₂* set, *test* set
3:     *Standard-Net* ← create a standard neural network  # *Net* is either FCN or CNN
4:     *Random-Net* ← best network obtained through random search
5:     *Evo-Single-Net* ← best network obtained through Evo-Single method
6:     *Evo-Triple-Net* ← best network obtained through Evo-Triple method
7:     *Coevo-Net* ← best network obtained through coevolution (Algorithm 1)
8:     **for** *net* in {Standard-Net, Random-Net, Evo-Single-Net, Evo-Triple-Net, Coevo-Net} **do**
9:         Train *net* on *train₁* + *train₂* sets for 30 epochs
10:         Test *net* on *test* set
11:         Record *net*'s *test* performance

**Figure 5: Overview of experimental coevolutionary setup. The coevolutionary algorithm uses the train$_1$ set to train a model and the train$_2$ set to evaluate the evolved model's fitness. The train$_1$ and train$_2$ sets are combined together and used as the training set for the best coevolved neural network. The final model is tested over the test set.**

CNN network was outperformed by the other methods in 15 out of the 16 groups of experiments, implying that automated AF search methods may increase model performance.

The fact that most wins involved 3 different AFs—either in Evo-Triple or in Coevo—affirms our surmise in Section 1: Combining different AFs improves network performance.

Table 7 presents the top-3 coevolved AF configurations for each of the four datasets. We note that some of the coevolved configurations are simpler and thus easier to use both for the feed-forward phase and the backpropagation phase. For brevity, we presented only the top-3 configurations, although there were many configurations that were composed only of simple evolved functions such as $x$, $2x$, and $x^2$.

## 6 CONCLUDING REMARKS

Despite the fact that deep-learning techniques are used for a wide range of applications, many hyperparameters must still be manually specified. The choice of AFs is a critical element that is often overlooked (i.e., standard AFs are used). We investigated the significance of AFs in learning in FCN and CNN models for image classification tasks. We introduced a coevolutionary algorithm to evolve new AFs and combine them beneficially, showing that our method performs well on four benchmark datasets.

We suggest a number of avenues for future research:

- Our focus herein was on classification tasks. Other deep-learning tasks can be examined too.

**Table 6: Experimental results. Each line represents 20 replicates. Values shown are mean accuracy percentages over test set, where the mean is computed over best per-replicate scores for a particular method. The winning method is marked in boldface. Leaky: Leaky ReLU; Arch: Architecture.**

| Dataset | Baseline | Arch | Standard | Random | Evo-Single | Evo-Triple | Coevo |
|---|---|---|---|---|---|---|---|
| MNIST | ReLU | FCN | 93.93 | 94.48 | 93.85 | 94.13 | **95.9** |
| | | CNN | 96.88 | 99.13 | 98.93 | **99.3** | 99.0 |
| | Leaky | FCN | 94.98 | 93.68 | 93.53 | **95.7** | 95.5 |
| | | CNN | 99.05 | 98.88 | 98.63 | 99.13 | **99.3** |
| KMNIST | ReLU | FCN | 80.3 | 80.58 | 80.88 | **81.33** | 80.15 |
| | | CNN | 59.58 | **95.15** | 94.78 | 94.68 | 94.8 |
| | Leaky | FCN | 78.33 | 80.15 | 81.23 | 80.45 | **81.68** |
| | | CNN | 94.7 | 93.98 | 93.15 | **95.38** | 94.23 |
| FashionMNIST | ReLU | FCN | 85.55 | 85.78 | 83.25 | 86.0 | **86.48** |
| | | CNN | 88.45 | 89.25 | 89.08 | 84.58 | **90.0** |
| | Leaky | FCN | 86.25 | 85.58 | 83.13 | **86.43** | 85.68 |
| | | CNN | **89.75** | 89.35 | 88.8 | 89.68 | 89.13 |
| USPS | ReLU | FCN | 90.52 | 91.27 | 91.39 | 84.91 | **91.4** |
| | | CNN | 89.1 | **96.13** | 93.77 | 95.39 | 96.01 |
| | Leaky | FCN | 92.52 | 91.77 | 92.27 | **92.89** | 92.39 |
| | | CNN | 96.38 | 95.39 | 94.39 | 96.01 | **96.88** |

**Table 7: Coevolution: Top-3 AF configurations per dataset, per architecture. Each AF is of the form $f(x) = ...$ with the $f(x)$ omitted for brevity. Each AF configuration is of the form {input-layer AF, hidden-layer AF, output-layer AF}.**

MNIST

FCN
$\{min(x, x^2) - max(x^2, \text{ReLU}(x)), \text{RReLU}(x), \text{ELU}(x)\}$
$\{x - max(x^2, \text{ReLU}(x)), x, x\}$
$\{x^2, \text{RReLU}(\text{RReLU}(\text{RReLU}(x))), 2x\}$

CNN
$\{x - \text{CELU}(x), \text{Hardswish}(\text{ReLU}(x)), \text{RReLU}(\text{ReLU}(\text{Hardshrink}(x)))\}$
$\{\text{Softshrink}(\text{Hardswish}(\text{Hardshrink}(x))), \text{RReLU}(x), \text{ReLU}(x)\}$
$\{\text{Hardtanh}(x) - \text{CELU}(x), \text{Hardswish}(\text{RReLU}(\text{ReLU}(x))), \tanh(x^2)\}$

KMNIST

FCN
$\{\text{ELU}(\text{Hardswish}(x)), \text{CELU}(x), \text{ReLU}(x)\}$
$\{max(0, \text{ELU}(x)), x, \text{ELU}(x)\}$
$\{\text{RReLU}(\text{ReLU}(\text{Hardshrink}(x))), min(\text{Softshrink}(\text{ELU}(x)), x), \text{ELU}(x)$

CNN
$\{x, \text{Softshrink}(\text{Hardtanh}(x)), \text{RReLU}(x) - \text{Hardswish}(x)\}$
$\{\text{Hardshrink}(x), \text{Softshrink}(\text{Hardshrink}(x)), \text{LeakyRelu}(\tanh(x^2))\}$
$\{min((\text{Hardshrink}(x) * x), \text{ELU}(x)), \text{Softshrink}(\text{Hardtanh}(x)), \text{Hardtanh}(\text{Hardswish}(x))\}$

FashionMNIST

FCN
$\{\text{RReLU}(x) \cdot x + x, \text{ELU}(x), 2x\}$
$\{x^2 + x, max(\text{ReLU}(x), \text{LeakyRelu}(x) + 2x) - (\text{LeakyRelu}(x) + 2x), 2x\}$
$\{\text{ELU}(max(\text{ReLU}(x), x)), \text{CELU}(\text{RReLU}(x)), \text{RReLU}(x)\}$

CNN
$\{\text{Hardswish}(x), \text{RReLU}(x), \text{ELU}(\text{LeakyRelu}(x) - x))\}$
$\{\text{Hardswish}(x), \text{Hardswish}(x), \text{Softshrink}(\text{Softshrink}(\text{CELU}(x)))\}$
$\{\text{Softshrink}(x), \text{Hardtanh}(\text{LeakyRelu}(\text{Hardswish}(x))), \text{Hardtanh}(max(\text{Hardtanh}(\text{Hardswish}(x)), 0))\}$

USPS

FCN
$\{\text{ELU}(\text{Hardswish}(x)), \text{CELU}(x), \text{ReLU}(x)\}$
$\{\text{Hardswish}(x), x, \text{ELU}(x)\}$
$\{\text{ReLU}(\text{CELU}(x)), \text{Hardswish}(\tanh(x)) + \text{RReLU}(x), \text{LeakyReLU}(\text{ELU}(x))\}$

CNN
$\{\text{RReLU}(\tanh(x)), \tanh(x), \text{LeakyReLU}(x)\}$
$\{\text{HardSwish}(\text{Softshrink}(x)), \text{CELU}(\text{ELU}(x) - min(\text{LeakyReLU}(x), x)), min(\text{Hardshrink}(x) - x, \text{LeakyReLU}(x))\}$
$\{\text{HardSwish}(\text{Softshrink}(x)), \text{ELU}(x) - min(\text{LeakyReLU}(x), x), \text{Softshrink}(min(\text{Softshrink}(x), \text{Hardshrink}(x) - x))\}$

- During DNN learning, the influence of AFs on computational complexity could be analyzed.
- Applying coevolutionary technique to other DNN hyperparameter searches.

- Inspecting other DNN topologies, such as recurrent neural networks (RNNs) and autoencoders (AEs).

# REFERENCES

[1] Forest Agostinelli, Matthew Hoffman, Peter Sadowski, and Pierre Baldi. 2014. Learning activation functions to improve deep neural networks. *arXiv preprint arXiv:1412.6830* (2014).

[2] Mina Basirat and Peter M Roth. 2018. The quest for the golden activation function. *arXiv preprint arXiv:1808.00783* (2018).

[3] James Bergstra and Yoshua Bengio. 2012. Random search for hyper-parameter optimization. *Journal of machine learning research* 13, 2 (2012).

[4] Garrett Bingham, William Macke, and Risto Miikkulainen. 2020. Evolutionary Optimization of Deep Learning Activation Functions. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference* (Cancún, Mexico) *(GECCO '20)*. Association for Computing Machinery, New York, NY, USA, 289–296.

[5] Garrett Bingham and Risto Miikkulainen. 2022. Discovering parametric activation functions. *Neural Networks* (2022).

[6] Thomas M Breuel. 2015. The effects of hyperparameters on SGD training of neural networks. *arXiv preprint arXiv:1508.02788* (2015).

[7] Yadong Ding, Yu Wu, Chengyue Huang, Siliang Tang, Fei Wu, Yi Yang, Wenwu Zhu, and Yueting Zhuang. 2022. NAP: Neural Architecture search with Pruning. *Neurocomputing* (2022). https://doi.org/10.1016/j.neucom.2021.12.002

[8] Giorgia Franchini, Valeria Ruggiero, Federica Porta, and Luca Zanni. 2023. Neural architecture search via standard machine learning methodologies. *Mathematics in Engineering* (2023). https://www.aimspress.com/article/doi/10.3934/mine.2023012?viewType=HTML

[9] Sepp Hochreiter. 1998. The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 6, 02 (1998), 107–116.

[10] Fernando Itano, Miguel Angelo de Abreu de Sousa, and Emilio Del-Moral-Hernandez. 2018. Extending MLP ANN hyper-parameters Optimization by using Genetic Algorithm. In *2018 International Joint Conference on Neural Networks (IJCNN)*. 1–8.

[11] Lu Lu, Yeonjong Shin, Yanhui Su, and George Em Karniadakis. 2019. Dying ReLU and initialization: Theory and numerical examples. *arXiv preprint arXiv:1903.06733* (2019).

[12] Sébastien Marcel and Yann Rodriguez. 2010. Torchvision the Machine-Vision Package of Torch. In *Proceedings of the 18th ACM International Conference on Multimedia*. Association for Computing Machinery, New York, NY, USA, 1485–1488.

[13] Risto Miikkulainen, Jason Liang, Elliot Meyerson, Aditya Rawal, Daniel Fink, Olivier Francon, Bala Raju, Hormoz Shahrzad, Arshak Navruzyan, Nigel Duffy, et al. 2019. Evolving deep neural networks. In *Artificial intelligence in the age of neural networks and brain computing*. Elsevier, 293–312.

[14] Julian Francis Miller and Simon L. Harding. 2008. Cartesian Genetic Programming. In *Proceedings of the 10th Annual Conference Companion on Genetic and Evolutionary Computation* (Atlanta, GA, USA) *(GECCO '08)*. Association for Computing Machinery, New York, NY, USA, 2701–2726.

[15] Julian F Miller and Stephen L Smith. 2006. Redundancy and computational efficiency in Cartesian genetic programming. *IEEE Transactions on Evolutionary Computation* 10, 2 (2006), 167–174.

[16] Julian F Miller, Peter Thomson, and Terence Fogarty. 1997. Designing electronic circuits using evolutionary algorithms. arithmetic circuits: A case study. *Genetic Algorithms and Evolution Strategies in Engineering and Computer Science* (1997), 105–131.

[17] Andrew Nader and Danielle Azar. 2021. Evolution of Activation Functions: An Empirical Investigation. *ACM Transactions on Evolutionary Learning and Optimization* 1, 2 (2021), 1–36.

[18] Chigozie Nwankpa, Winifred Ijomah, Anthony Gachagan, and Stephen Marshall. 2018. Activation functions: Comparison of trends in practice and research for deep learning. *arXiv preprint arXiv:1811.03378* (2018).

[19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. PyTorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019), 8026–8037.

[20] Carlos Andrés Pena-Reyes and Moshe Sipper. 2001. Fuzzy CoCo: A cooperative-coevolutionary approach to fuzzy modeling. *IEEE Transactions on Fuzzy Systems* 9, 5 (2001), 727–737.

[21] Markus Quade. 2020. cartesian. https://github.com/Ohjeah/cartesian.

[22] Snehanshu Saha, Nithin Nagaraj, Archana Mathur, and Rahul Yedida. 2019. Evolution of novel activation functions in neural network training with applications to classification of exoplanets. *arXiv preprint arXiv:1906.01975* (2019).

[23] Sagar Sharma and Simone Sharma. 2017. Activation functions in neural networks. *Towards Data Science* 6, 12 (2017), 310–316.

[24] Moshe Sipper. 2021. Neural Networks with À La Carte Selection of Activation Functions. *SN Computer Science* 2, 6 (2021), 1–9.

[25] Moshe Sipper, Jason H Moore, and Ryan J Urbanowicz. 2019. Solution and fitness evolution (SAFE): Coevolving solutions and their objective functions. In *European Conference on Genetic Programming*. Springer, 146–161.

[26] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. 2012. Practical Bayesian optimization of machine learning algorithms. *arXiv preprint arXiv:1206.2944* (2012).

[27] Kenneth O Stanley and Risto Miikkulainen. 2002. Evolving neural networks through augmenting topologies. *Evolutionary computation* 10, 2 (2002), 99–127.

[28] Yanan Sun, Bing Xue, Mengjie Zhang, and Gary G Yen. 2019. Evolving deep convolutional neural networks for image classification. *IEEE Transactions on Evolutionary Computation* 24, 2 (2019), 394–407.

[29] Meng-Ting Wu, Hung-I Lin, and Chun-Wei Tsai. 2022. A Training-Free Genetic Neural Architecture Search. In *Proceedings of the 2021 ACM International Conference on Intelligent Computing and Its Emerging Applications* (2022-01-01) *(ACM ICEA '21)*. Association for Computing Machinery, Jinan, China, 65–70. https://doi.org/10.1145/3491396.3506510

[30] Haoyu Zhang, Yaochu Jin, Yaochu Jin, and Kuangrong Hao. 2022. Evolutionary Search for Complete Neural Network Architectures with Partial Weight Sharing. *IEEE Transactions on Evolutionary Computation* (2022), 1–1. https://doi.org/10.1109/TEVC.2022.3140855

[31] Xiaohu Zhang, Yuexian Zou, and Wei Shi. 2017. Dilated convolution neural network with LeakyReLU for environmental sound classification. In *2017 22nd International Conference on Digital Signal Processing (DSP)*. IEEE, 1–5.

[32] Yu-Dong Zhang, Chichun Pan, Junding Sun, and Chaosheng Tang. 2018. Multiple sclerosis identification by convolutional neural network with dropout and parametric ReLU. *Journal of computational science* 28 (2018), 1–10.