

## SYNTHETIC CFD ESTIMATION FOR BLACKHAWK AIRFOIL DRAG COEFFICIENT

James Ross\*

James.E.Ross@erdc.dren.mil

Robert Haehnel†

Robert.B.Haehnel@erdc.dren.mil

Luke Allen†

Luke.D.Allen@erdc.dren.mil

Stephen Wang†

swang.2420@gmail.com

Brianna Thompson\*

Brianna.D.Thompson@erdc.dren.mil

Yi Zhang\*

Yi.Zhang@erdc.dren.mil

### **ABSTRACT.**

- TODO: Questions 1. We only really study drag in the paper. maybe title change to only relate to drag and mention lift and moment as areas we can apply our method to for future work.  
2. Do we do well in (hyperparameter subsection), "key areas of concern, and in this example, at low drag coefficient values."  
3. Is the c81gen code completely custom or is it a standard that was modified? adjust abstract description.

Through understanding the effects of the modifications to the geometry of a rotorcraft blade airfoil, we will be able to better model airfoil performance. To analyze the performance of different airfoils, a morphing tool, Parfoil, is used to produce an airfoil geometry. The performance of the airfoil geometry is calculated using 2D a Computational Fluid Dynamics (CFD) model in the C81gen computer code. The effects of airfoil modification on rotor blade performance is then assessed using a rotorcraft comprehensive analysis (CA) tool and successively refined with a numerical optimizer. A critical time constraint in the process of approximating airfoil performance is the calculation of the airfoil performance tables that the CA tool requires. The time constraint occurs as a result of the need to use computationally expensive CFD models for generating the airfoil performance tables (tabulation of airfoil force coefficients at specific relative air speeds and blade angles of attack) required by the CA tool for determining vehicle performance metrics. As an alternative to directly generating these "air tables" for each evaluation point with CFD, use of a machine learning (ML) surrogate model is explored in this work. This method would decrease analysis time as thousands of air tables are needed to determine the optimal design. A base ML surrogate model was created to determine the feasibility of using a ML surrogate model to estimate the performance of an aircraft airfoil. The base ML surrogate model showed a respectable  $R^2$  output for error from the loss function during training, but upon closer inspection during the validation step of training the initial network, many predictions still showed an unacceptable magnitude of error. The objective of this work is to improve the ML surrogate model's accuracy so that it becomes a feasible alternative for

---

\*U.S. Army Engineer Research and Development Center, Vicksburg, Mississippi, U.S.A.

†U.S. Army Engineer Research and Development Center, Hanover, New Hampshire, U.S.A.

estimating the airfoil tables that would ordinarily be produced by running CFD simulations. Because the computational time requirements of running CFD models for estimating the rotorcraft airfoil performance are restrictive, having a ML surrogate model would allow for estimating a large number of potential airfoil designs in a much shorter timeframe. A ML surrogate model would also allow for a more efficient search of possible designs leading to a more robust final design. This final design could potentially reduce the number of flight tests that were previously required to find a similar design. The results defined in this work are an analysis of the effects of modifying the hyperparameters of the base ML surrogate model to improve accuracy seeking results comparable to CFD combined with linear interpolation methods.

**Keywords:** Machine Learning, Rotorcraft, Optimization

## 1 Introduction

Rotorcraft designers often use comprehensive analysis codes, such as RCAS [30, 16] and CAMRAD II [17], to evaluate and predict a wide range of characteristics and performance information for their models (e.g., rotor trim, aerodynamic loads, vibrations). These codes typically do not include fluid flow solvers and thus rely on tabulated data that define the airfoil lift, drag, and pitching moment coefficients ( $c_l$ ,  $c_d$ , and  $c_m$ , respectively) for each airfoil throughout the operational regime. The tables may use actual flight test data or data that is calculated via other means, such as using computational fluid dynamics (CFD). CFD calculations, however, are computationally expensive to produce, making them intractable for applications such as design optimization where the number of variations may be very large. Researchers have attempted to mitigate this issue by employing various modeling techniques to reduce the overall computation time. One such effort by Allen et al. [2] used C81Gen [28, 22] to generate a database of airfoil tables for a range of designs and then used linear interpolation to estimate the values of  $c_l$ ,  $c_d$ , and  $c_m$  for any airfoil geometry within the design space. Their work recognized the limitations of the linear interpolation scheme and proposed that machine learning could be used as an alternative. The current effort was designed to explore this idea in detail and provide a framework for how such a problem can be approached.

We began this effort with a base hand configured neural network which resulted in a network that suffered from overfitting. The idea was then to apply hyperparameter optimization to the hand designed neural network to make an attempt at finding a more accurate neural network design that did not suffer from overfitting (see Section 3.2). Furthermore, we note that the rotor performance metric being used by Allen et al. (the rotor power coefficient,  $C_p$ ) is most sensitive to the airfoil drag. Therefore,  $c_d$  has been the focus of testing and validation to this point, though the same methodology can be readily applied to  $c_l$  and  $c_m$ . In this section we will introduce the concepts and provide a small description of Linear Interpolation 1.1, Machine learning 1.2, Original CFD Surrogate Model Analysis 1.3, and Hyperparameter Optimization 1.4.

## 1.1 Linear Interpolation

Linear interpolation is commonly used for estimating new data points in between connecting pairs of known points determined from experiment or using CFD [23]. This operation is simple and fast, so it would significantly cut down on computational time when a large number of air tables are needed. Relative to the number of points produced, time required to perform CFD runs increases linearly while time required for interpolation remains constant. The graph in Figure 1 shows a time comparison of a pure CFD run against the interpolation method based on the number of air tables that need to be generated. The speedup of the interpolation method starts exceeding using only CFD by an order of magnitude with 10,000 tables, and by two orders of magnitude with 100,000 tables. However, the benefit of interpolation disappears with 1,000 tables or less due to a base number of CPU hours required to generate the database for interpolation; in the test case, 304 tables were created with CFD to base points to facilitate linear interpolation.

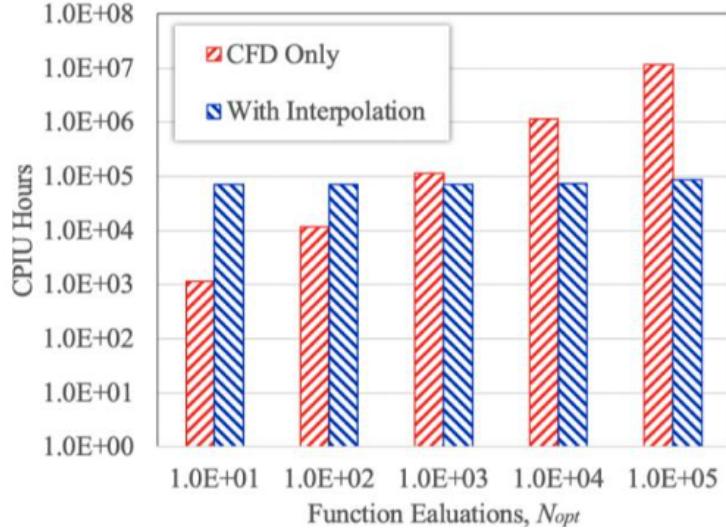


Figure 1: CFD Compare With Interpolation

Though linear interpolation is computationally efficient, reasonably accurate (provided the grid spacing is not too large) and well behaved, linear interpolation to parameter space greater than 3 starts to become unwieldy [27, 3]. Furthermore, linear interpolation schemes are most efficient on a regularized grid structure. However, it is often advantageous to sample the parameter space in a non-uniform way, concentrating more points in areas of most interest while sparsely sampling regions that are perceived to be less fruitful. The efficiency of linear interpolation breaks down when nearest neighbor searches consume significant amounts of computational time on large, irregularly sampled parameter spaces. Alternative surrogate models, such as machine learning, can be

more efficient on irregularly sampled parameter spaces.

## 1.2 Machine Learning

We are experiencing a revolution in the field of machine learning, which has been enabling advances across a wide range of scientific and engineering areas [36, 29]. Machine learning is focused on the development of algorithms with the ability to build models from obtained data without explicit use of mathematical models [35]. Applying machine learning methods to CFD is focused on the possibility of increasing the speed of highly accurate simulations, developing turbulence models with higher accuracy, and producing advanced models, which is beyond what can be achieved using classical approaches [12]. Some researchers have proposed ML methods for CFD and obtained good results, such as the turbomachinery design, turbulence modeling, and heat-transfer for aerodynamic optimization [38, 12, 11]. Among the machine learning methods is the use of a neural network, a popular supervised learning algorithm for modeling complicated systems [14, 4, 39]. Neural networks are based on a collection of connected nodes called neurons, which are arranged in parallel layers that are fully interconnected by weighted connection. The most important characteristic of a neural network is learning from data [10].

## 1.3 Original CFD Surrogate Model Analysis

A ML surrogate model, more specifically, a regression-based neural network (RNN), receives a dataset as input. Within the RNN there are many parameters which process the dataset, and these parameters are known as neurons. Each of these neurons is tuned and collectively determine a function for predicting data values. Neurons are positioned in a layered data structure. The first layer is known as the input layer, and the only job it has is to accept the incoming data source. Then, there are one or more layers known as hidden layers which are used to introduce non-linearity into the network. After hidden layers, there is a single layer known as the output layer, which provides the prediction of the network. Each layer contains an activation function within its neurons, which determines how the neurons are used in the calculation of the network's prediction. In summary, the configuration of the network's neurons, the number of layers, and the activation function within each neuron are used to process the data and provide a value to the user.

Our base surrogate ML model was designed manually with 16 neurons in the input layer, three hidden layers, and a single-neuron output layer. We soon realized this model was not sufficient on the dataset we provided; 31.8% of the data points had more than 10% error. The original model also performed better on the training dataset than on the test dataset which is commonly referred to as overfitting (see Section 3.2). This poor performance on test data but great performance on training data is due to the model learning too much about the training data.

## 1.4 Hyperparameter Optimization

Since the original surrogate model suffered from issues with being able to accurately predict because of overfitting, we looked to improve on the base model through the use of hyperparameter optimization. Hyperparameter optimization is the process of sweeping through values for model parameters to find a set of parameters that produce an answer within an acceptable error range. Neural networks are not always accurate when making predictions, but could be described as producing an expert-level human like response to a question. Therefore, the goal of our use of hyperparameter optimization was to find a network that not only gave good quality of fit metrics (low RMSE and high  $R^2$  values for the entire model), but was accurate in key areas of concern, and in this example, at low drag coefficient values. Our search will seek to optimize hyperparameters important to the network design with conditions listed below:

- Number of hidden layers/neurons in each layer: How many parameters does the model have, and how deep is it? (Section 2.4)
- Hidden layer activation: How is non-linearity introduced in the model? (Section 2.5)
- Output layer activation: How should the model map the final result? (Section 2.6)
- Optimizer: What learning optimizing algorithm should the model use? (Section 2.7)
- Learning rate: What is the response level of the model to the estimated error of the output? (Section 2.8)
- Epochs: How many times should the model be trained before it starts overfitting? (Section 2.9)
- Loss function: How is model performance quantified? (Section 2.10)

## 2 Neural Network Components and Selections

In a prior work, Allen et al. (2021) used linear interpolation to estimate the air tables from the database generated from CFD [2]. While interpolation is an efficient calculation, its accuracy depends heavily on how far apart each data point is from its paired point. Using a finely resolved regular grid of CFD runs would produce closely-spaced points, but is not possible due to time constraints. As such, a practical CFD run would leave us with wider gaps between data points, potentially leading to greater inaccuracies in interpolated points far from a CFD computed point. Given this limitation, our objective shifted towards creating a ML surrogate model that can accurately reproduce points in the CFD created database, as well as accurately interpolate points for locations not contained in the database. Another requirement of the the ML model would be to compute required air tables at a rate several orders of magnitude faster than CFD. By using a ML model to provide a reasonably accurate estimate for the performance table

values produced by a CFD model, we sought to significantly reduce the computational effort required for generating the airfoil performance tables without the need to run resource intensive CFD models. We were also looking for an improvement in accuracy of expected values in between points produced from CFD runs over the expected values produced by the linear interpolation method. In the following section we are going to describe the database, neural network libraries, hyperparameter optimization libraries, and the neural network components.

## 2.1 Dataset

«< TODO: ask Luke and Bob if compilation is correct or does lift, drag, and moment have their own values. if they do I would like to rephrase this to say drag has been used but lift and moment will be reviewed in future work »>

The dataset is composed of air tables which are a compilation of lift, drag, and moment coefficients as a function of free stream Mach number ( $M$ ) and angle of attack (AoA). Once these values for the air tables are generated using CFD, those values are provided to the ML model for training. In total, we generated 304 tables, which contained 190,000 training points. Ninety percent of that data is allocated to the training set, which the regression-based neural network (RNN) uses to determine a non-linear function for predicting new air tables for new airfoil shapes. The remaining ten percent of the values within the air tables become the testing set, which is used for evaluating expected performance of the model on unseen data.

## 2.2 Neural Network Libraries

The Keras library was used for creating the neural network used in this work. Keras, which provides easy to use commands to interface with Google’s TensorFlow library, provides a clean interface to build regression-based neural network (RNN) models [19]. Though Facebook’s PyTorch is another popular choice, Keras was used since the team had more experience with this tool.

## 2.3 Hyperparameter Optimization Libraries

Because a regression-based neural network (RNN) has endless configurations, it is highly unlikely to find the optimal RNN by hand. To more efficiently explore the search space, we decided to use hyperparameter tuning libraries. These tools can automate testing of different RNNs based on hyperparameter conditions that we predetermine. Two such libraries that stood out were Optuna and KerasTuner. Both have a simple syntax to set up the optimization process and include state-of-the-art search algorithms. KerasTuner is integrated within Keras, so this library was the natural starting point. However, we found some of the functionality within Optuna to be helpful for this this work. For example, Optuna allows a user to store trials in a SQL database, allowing for organized viewing of past trials and checkpoints to resume optimization. On the other hand, the KerasTuner storage format was less accessible and used several times more storage when

tested for this effort. Pruning is another feature of Optuna that was helpful in performing hyperparameter optimization for this work. Pruning ends unpromising trials early based on a certain condition (e.g. if its intermediate value is worse than the median of past trials), which cuts down on optimization time. Because of efficient storage and trial execution, Optuna was chosen for hyperparameter tuning [26, 32]. Optuna has several search algorithms for determining hyperparameters. For hyperparameters with a limited search space (e.g. optimizer, activation function, and loss function), we used grid search to perform an exhaustive evaluation. For hyperparameters with an potentially large search space (e.g. number of layers and neurons in each layer, learning rate, epochs), we used Bayesian search, which employs probabilistic measures to narrow in on optimal regions.

## 2.4 Number of Hidden Layers and Neurons per Layer

For a regression-based neural network (RNN), the optimal number of hidden layers and neurons per hidden layer largely depends on the problem. Not having enough hidden layers/neurons can result in underfitting, which means the model learns too little about the training data to adequately make predictions. On the other hand, having too many neurons in the hidden layers may cause overfitting, meaning the model learns patterns in the training data so well that it is unable to generalize to unseen testing data. Most problems can see reasonable results with 1-5 hidden layers, depending on the trade-off one wants between accuracy and computation time [34]. Within the hidden layers, there are significant possibilities for the number of neurons inside them, so we let Optuna automate combinations based on conditions we specified. The results can be seen in Section 3.3.1.

## 2.5 Hidden Layer Activation

ReLU has become the standard for hidden layer activation due to the function’s simplicity and ability to commonly outperform other activations like Sigmoid and Tanh on image classification benchmarks. Despite its many strengths, new activation functions have been proposed to improve upon ReLU, including Google’s Swish and Misra’s Mish. The base model originally used Mish as the hidden layer activation. Mish has properties that would theoretically make it more ideal than alternative activation functions, like allowing for negative gradients, having a continuous first derivative, and being unbounded above and bounded below [25]. For the original base model Swish was also tried, but the model training loss exploded to infinity, so it was not used in the final network design. To determine whether to use Mish or ReLU, we tested each on the same model for 10 trials, with 25 epochs per trial, and the test results can be seen in Section 3.3.2. We chose ReLU for the hidden layer activation function as it produced the lower average validation loss.

## 2.6 Output Layer Activation

The base model used softmax activation, but we soon realized that this was not the most appropriate choice. Softmax is used often in multi-class classification problems because it outputs a probability distribution. However, we have a regression problem, so we instead used a linear activation for the output layer.

## 2.7 Optimizer

Initially, the Adam optimizer was used for training the regression-based neural network (RNN). Adam improves upon stochastic gradient descent through use of a dynamic learning rate, which allows for more effective convergence and less chance of being stuck in local minima. The reliability of Adam makes it a solid choice for many practical deep learning problems. Kingma and Ba showed that Adam outperformed four other optimizers on MNIST and CIFAR-10, two benchmark image classification datasets [20]. Although we have a regression problem and not a classification problem, Adam still makes a good starting point given that it is known to work well on RNNs.

It was initially suspected that the Adam optimizer would outperform other optimizers for our particular problem. This suspicion that Adam would outperform other optimizers was related to Adam generally performing very well for RNNs. Eventually we let Optuna do a grid search on eight different optimizers (SGD, RMSprop, Adam, Adadelta, Adagrad, Adamax, Nadam, Ftrl) built into Keras to see if there was one that might outperform Adam. Each optimizer was tested on the same model for 25 epochs. Then, the top three optimizers (Adam, Adamax, Nadam) were tested further with 50 and 100 epochs averaged over three trials. The outcome of the optimizers test can be observed in Section 3.3.3.

## 2.8 Learning Rate

The default learning rate (LR) for Adamax is 0.002, as suggested by Kingma and Ba. However, the majority of Keras's built-in optimizers have a default of 0.001. To decide the LR, we let Optuna perform a grid search in the range from 0.0001 to 0.002 with step size of 0.0001. The top 10 LRs for 10 epochs were further tested for 25 epochs. Then, the top 4 from the 25 epochs were compared for another 100 epochs. The results from the test performed on LR can be seen in Section 3.3.4.

## 2.9 Number of Epochs

Throughout our testing, we gradually increased the number of epochs in training as we found more promising configurations. Later on, we employed Keras's early stopping feature, as it helped us to thoroughly train the model while helping to prevent overfitting. Early stopping takes a parameter called patience, which stops model training after a certain number of epochs if the validation loss does not improve. Patience is commonly set somewhere between 1 and 100, so in our case we arbitrarily chose 30. It is possible

that there is a better patience value for our network, however more testing for a better patience value will be left to future work. The results for the testing different options for the number of epochs can be seen in Section 3.3.5.

## 2.10 Loss Function

Part of the hyperparameter optimization process led us to look at different loss functions as a possibility for improving the predictions provided by our model. We performed several studies which focused on the use of the Mean Absolute Percent Error (MAPE) loss function. MAPE produced good results for our predictions except in the data range close to zero. Hyperparameter optimization led to the examination of the Huber and LogCosh loss functions. The results of test performed while searching for an optimal loss function can be observed in Section 3.3.6.

# 3 Results

In this section we will show results for the base model and discuss methods for combating overfitting. To improve the base model's performance we investigated a number of options for the number of hidden layers and neurons per layer, hidden layer activation, optimizer, learning rate, number of epochs, and loss function through hyperparameter optimization. We will show the test results from the experiments carried out to determine the optimal hyperparameters for model improvement. Then, we will present the steps we took to reach our final model and provide a comparison of the results of the base model and final model. Finally, we will introduce ensemble methods created using models produced during hyperparameter optimization, describe the tests we performed using ensemble methods, and show the results from these tests.

## 3.1 Base Model

The graph in Figure 2 shows the percent error of the base model's predictions of drag coefficient on 500 data points, with a logarithmic y-axis. Points are colored based on the angle of attack of the airfoil. We see a large data cluster around .15 to .25 actual value with a general expected error of an order of magnitude above the actual value. There were also a significant number of points that were up to two orders of magnitude off from the actual value.

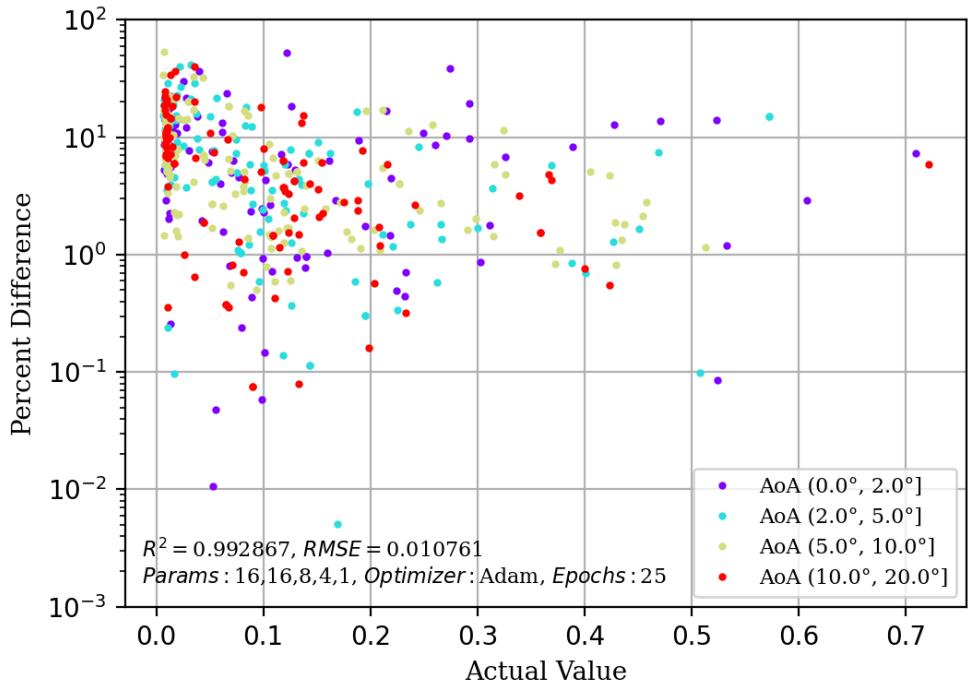


Figure 2: Drag Original 25a

The first predictions of the best-performing regression-based neural network (RNN) for estimating airfoil performance had unacceptable magnitude of error and produced inaccurate results. In the next sections, we compare the initial design and each of its modifications that led to the architecture of the final ML model. The choices to follow were informed by tests and visualizations generated in Optuna and TensorBoard.

### 3.2 Overfitting

One popular method to prevent overfitting is to introduce a dropout layer. Dropout randomly turns off neurons inside a neural network to avoid unintentional dependencies between neurons [31]. However, we found that dropout did not work well on our ML model. Instead, we used TensorBoard, a visualization tool integrated within TensorFlow, to verify that new models we found were not overfitting [33]. In the imaginary example below, a model with loss curves of this behavior should be trained for not more than 100 epochs based upon a convergence and the divergence of the curves shown in Figure 3.

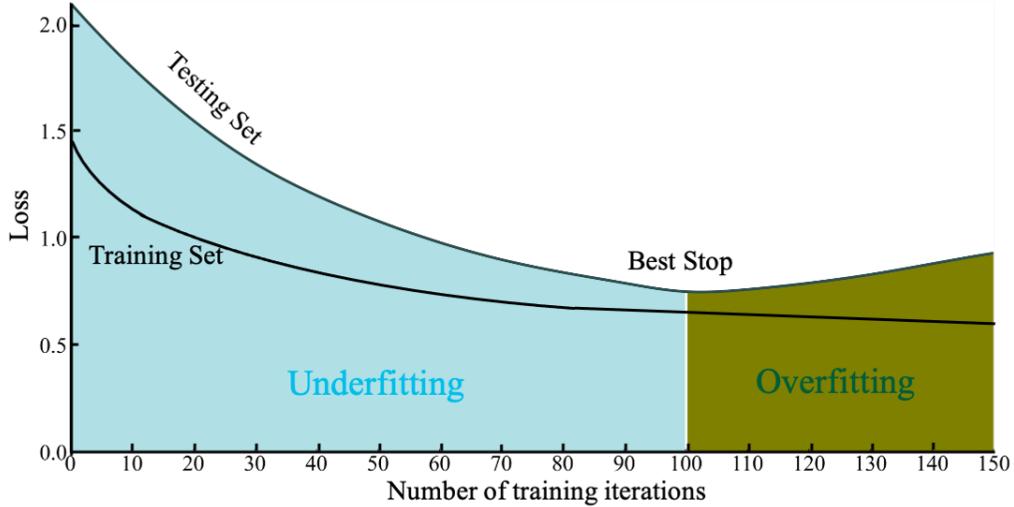


Figure 3: Overfitting Vs Underfitting

When we increased the number of epochs, number of hidden layers, or number of neurons per layer, we made sure that the model’s loss curves in TensorBoard did not exhibit the behavior in the green region above. More specifically, we were looking to identify overfitting in the case where the training data continued to improve but the results in predictions from the testing set started to become more inaccurate.

### 3.3 Neural Network Architecture Results

In this section we present the results from the hyperparameter search using Optuna and discuss the final choice for each hyperparameter.

#### 3.3.1 Number of Hidden Layers and Neurons per Layer

In an effort to improve our overall model results, we decided to implement the use of Optuna. Optuna is a great tool for sweeping through parameters and looking for better network conditions. Within Optuna we sent in value ranges for the number of hidden layers and the number of neurons in each layer. The following are the initial conditions we passed into Optuna for hyperparameter optimization:

Number of hidden layers: 3-7

Input layer size: 14-64 neurons

Hidden layers size: 0.5x-1x of the previous layer (e.g. if the input layer has 60 neurons, the next layer would have between 30 to 60 neurons)

The parallel coordinate plot in Figure 4 shows the individual trial parameters that Optuna found with the initial set of conditions. The parameters in Figure 4 can be described

as follows:

Objective value	Validation loss
n_layers	Number of hidden layers
n_units_input	Number of neurons in input layer
n_units_layer_0	Number of neurons in hidden layer 1
n_units_layer_1	Number of neurons in hidden layer 2

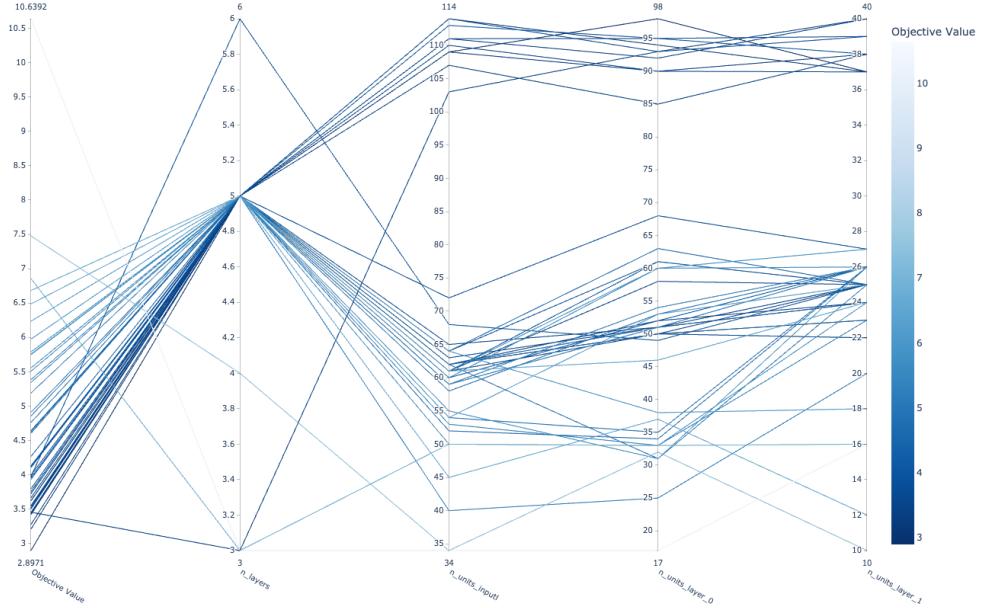


Figure 4: MAPE Parallel Coordinate Plot for the MAPE Loss Function

In total, 760 trials were run, 709 trials were pruned, and 51 trials were completed. We see in Figure 4 that Optuna chose 5 hidden layers as the optimal number with input layer sizes in the ranges of [40, 72] and [107, 114] neurons. For the hidden layer sizes, there was no clear pattern for the first one or two hidden layers, but deeper into the network the layers tended to be about half the size of the previous layer.

After examining the results from the initial conditions, we decided to further increase the complexity of the model since it was still training very fast with 5 hidden layers. With the following conditions, Optuna was used to decide the configuration that would become the final model:

Hidden Layers	6-12
Input Layer Size	100-400 neurons
Hidden Layer Size	2-400 neurons

The parallel coordinate plot in Figure 5 shows the individual trial parameters that Optuna found with the conditions for the final model. 389 trials were run, 371 were pruned,

and 18 were completed. Increasing the number of hidden layers allowed Optuna to find an input layer in the range of [329, 399] neurons, with 6 hidden layers being optimal. While 7 layer configurations were completed more often and had quite a few good objective values, it also produced several of the worst objective values. More than 7 layers likely results in diminishing returns and overfitting, with only three such completed trials. The hidden layers still decreased in size as the network got deeper, despite the condition allowing a layer to have more neurons than the previous.

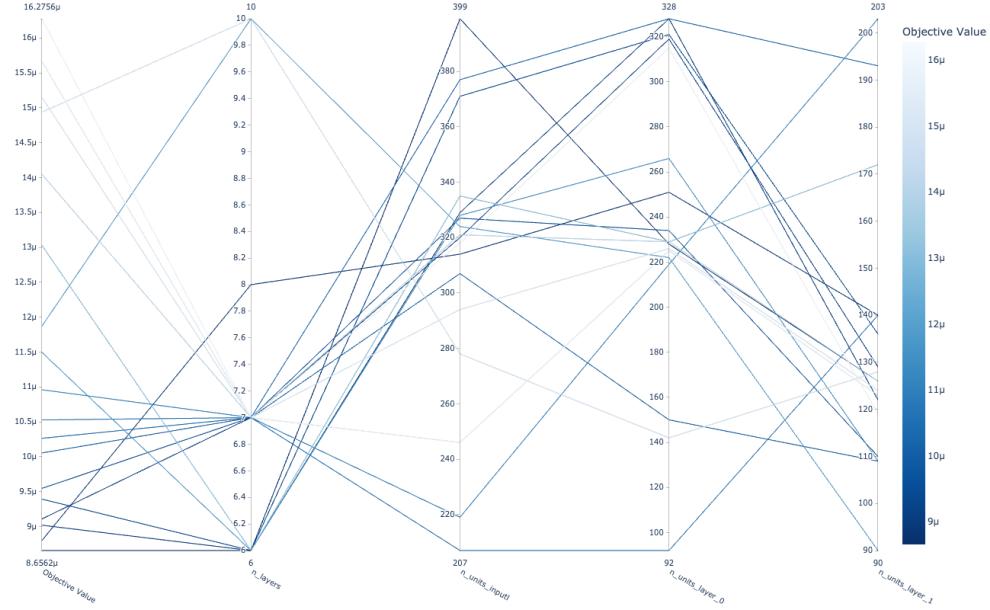


Figure 5: Parallel Coordinate Plot for the Conditions for the Final Model

### 3.3.2 Hidden Layer Activation: Mish vs. ReLU

The results of the hyperparameter search for the hidden layer activation function are presented in Table 1, showing base model runs comparing Mish and ReLU (highlighted). The average validation loss for ReLU was 2.755%, compared to 3.097% for Mish. While 4 out of the top 10 runs were Mish, the activation also produced the bottom 5 runs, showing that Mish is less consistent than ReLU.

Rank	Activation	Val Loss (25E)	Rank	Activation	Val Loss (25E)
1	ReLU	2.4988%	11	ReLU	2.8479%
2	Mish	2.50%	12	ReLU	2.8751%
3	ReLU	2.5132%	13	Mish	2.99%
4	ReLU	2.5478%	14	ReLU	3.0354%
5	Mish	2.61%	15	ReLU	3.0434%
6	Mish	2.62%	16	Mish	3.28%
7	ReLU	2.6912%	17	Mish	3.36%
8	ReLU	2.7399%	18	Mish	3.59%
9	Mish	2.75%	19	Mish	3.62%
10	ReLU	2.7609%	20	Mish	3.65%

Table 1: Hidden Layer Activation Test

In practice, Mish outperforms ReLU in some cases on the benchmark MNIST digits dataset, but the results came from training networks with 15-25 layers and varying batch sizes. ML problems do not often need so many layers, and in our case, we did not need to split our data into batches. Mish would be more appealing had there been favorable comparisons on smaller network sizes. We decided to switch back to ReLU given its reliability and our training results shown in Table 1.

### 3.3.3 Optimizer: Adam vs. Nadam vs. Adamax

Table 2 shows the results of testing the same model on different optimizers, which covers all of those built into Keras. We see that Adamax, Adam, and Nadam came out on top for 25 epochs.

Optimizer	Val Loss (25E)	Val Loss (50E)	Val Loss (100E)	Time (100E)
Adamax	1.02%	0.85%	0.68%	10.3 min.
Adam	1.49%	1.11%	0.87%	10.5 min.
Nadam	1.95%	1.09%	0.82%	14.5 min.
RMSProp	2.44%	-	-	-
Adagrad	4.24%	-	-	-
Ftrl	5.53%	-	-	-
Adadelta	11.70%	-	-	-
SGD	NaN	-	-	-

Table 2: Optimizers Test

The performance of Nadam and Adamax are not surprising, as they have been closely comparable to Adam in several studies. Kandel et al. tested these three optimizers among others in histopathology image analysis and found that each excelled depending on the type of network architecture used (Adam - VGG16, DenseNet; Nadam - ResNet; Adamax - InceptionV3). However, 3 of the 4 tests did not show statistically significant results

between the best and worst performance of the three optimizers, with accuracy differences of 0.15%, 0.66%, 0.81%, and 1.49%. The 1.49% difference was found on InceptionV3 network using Adamax, which may suggest its superiority in certain situations [18]. Another study conducted by Dogo et al. tested these optimizers on three randomly-selected Kaggle image datasets. Nadam came out on top for two of the three datasets, while the remaining one went to Adam. The study also compared convergence times between the optimizers, and found that Adamax converged the fastest, followed by Adam and then Nadam [9]. Ultimately, there seems to not be a clear choice between Adam, Nadam, and Adamax; it more so comes down to what works best for the problem at hand. For our study, we found that Adamax produced the best results with fast convergence (as seen in Table 2).

### 3.3.4 Learning Rate

Table 3 shows the results of testing different learning rates (LR) on the same model. A LR of 0.0006 produced the best result for 10, 25, and 100 epochs. For 100 epochs though, the validation losses for LRs of 0.0009, 0.0019, and 0.002 trailed close behind. Still, at 100 epochs 0.0006 had the best value percent loss at 0.856%. Generally, a lower LR achieves a smoother convergence at the expense of more epochs. Here however, a LR of 0.0006 is still able to converge in a reasonable amount of time, so the time tradeoff is not as relevant as would normally be expected [40].

LR	Val Loss (10E)	Val Loss (25E)	Val Loss (100E)
0.0009	1.12%	0.98%	0.89%
0.0014	1.21%	1.24%	-
0.002	1.28%	0.92%	0.88%
0.0019	1.30%	0.96%	0.91%
0.0006	1.311%	0.847%	0.856%
0.0007	1.33%	1.12%	-
0.0012	1.34%	1.11%	-
0.0003	1.41%	-	-
0.0008	1.47%	-	-
0.0005	1.51%	-	-

Table 3: Learning Rate Test

### 3.3.5 Number of Epochs

Figure 6 shows the TensorBoard graph for the best model run for 1000 epochs. We see that in the first 200 epochs, the validation loss curve closely follows the training loss curve. Beyond 200 epochs, the validation curve starts to plateau and separate from the training curve, suggesting signs of overfitting. Keras's early stopping callback on average stopped the model at around 200 epochs, so that is the number we used to train the best model.

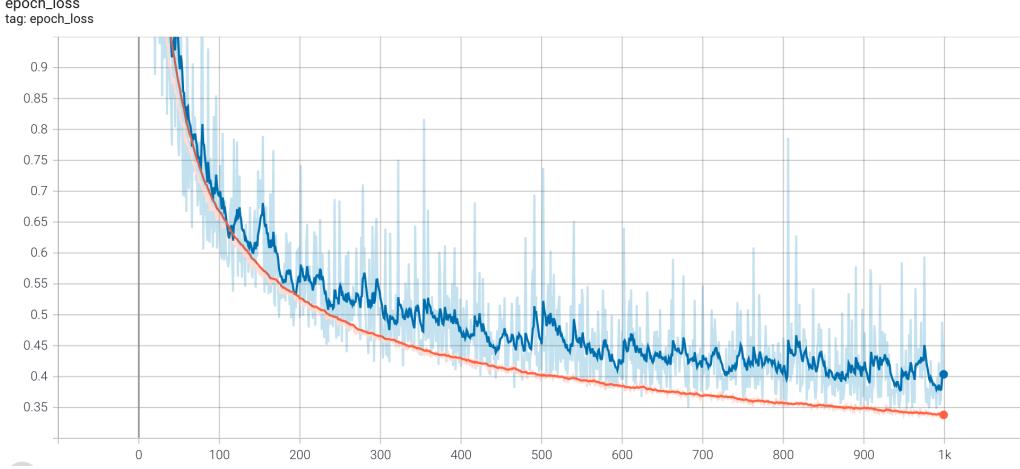


Figure 6: Number of Epochs

### 3.3.6 Loss Functions

Although this study focused mainly on drag, leaving the analysis of lift and moment to future work, we did perform some basic tests of our model against lift and moment. During our look at the loss function, we noticed that there was a limitation of MAPE on data points close to 0 when performing initial training of the lift and moment models; it seemed as though the weights were failing to update properly with MAPE in these models (Figure 7). We considered Huber and Logcosh as alternative loss functions, and the results for Logcosh are shown in Figure 8. The use of Logcosh repaired the issues with training lift and moment seen when using the MAPE function and produced more reasonable looking plots (Figure 8). In future work for lift and moment models, Logcosh will be preferred over Huber. Though they share the same advantages, the Huber loss function would require tuning an additional hyperparameter, which could be costly to determine.

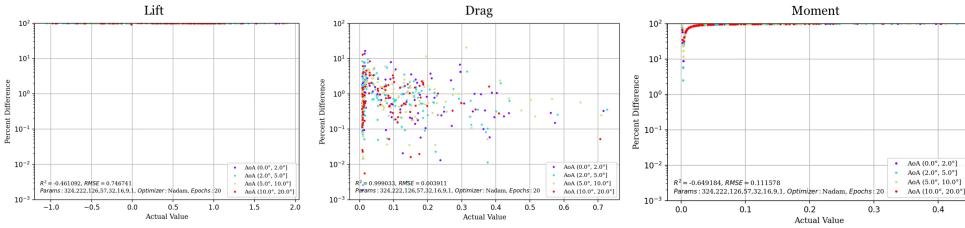


Figure 7: MAPE Loss Function Test Results

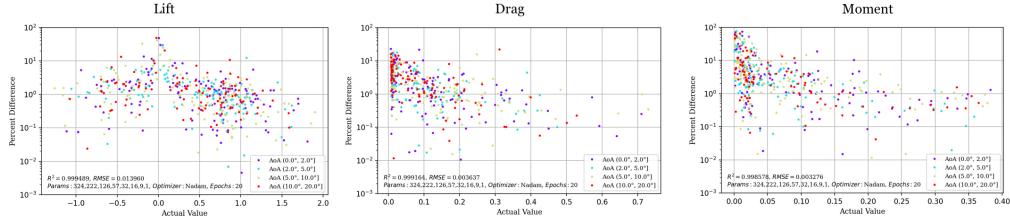


Figure 8: Initial Logcosh Loss Function Test Results

### 3.4 Summary of Steps to Final Model



«< TODO: need clarification from stephen what table 4&5 are supposed to show. it is unclear. This section will need a clarification pass to better describe the figures after talking with stephen bob and luke»>

Here we will look at a summary of the steps taken to get from the starting point of linear interpolation, transitioning to building the base model using a neural network, followed by utilizing hyperparameter optimization. We will then take a short look at the best model located through this process before ensembling methods were applied.

#### 3.4.1 Linear Interpolation

<b>mu</b>	<b>0.0</b>	<b>0.05</b>	<b>0.1</b>	<b>0.15</b>	<b>0.24</b>	<b>0.3</b>	<b>0.35</b>	<b>0.37</b>	<b>0.4</b>
count	102	102	102	102	102	102	102	102	102
mean	0.8534	0.8605	1.1152	1.074594	1.147324	1.433962	1.272401	1.218535	1.107003
std	0.2493	0.2043	0.3042	0.722038	3.281998	0.703469	0.543612	0.68337	0.70893
min	0.2463	0.44157	0.540728	-0.577864	-4.90635	-1.81702	-0.53454	-1.96081	-0.61557
25%	0.6648	0.705023	0.885208	0.812251	0.73259	1.104198	0.965715	0.823362	0.707655
50%	0.8294	0.840901	1.045421	1.170881	1.249883	1.376993	1.240295	1.215328	1.045802
75%	1.0240	0.980409	1.326864	1.560228	2.281964	1.877479	1.654725	1.629042	1.393702
max	1.471513	1.420221	2.08377	2.501025	8.219547	2.744838	2.559548	2.948367	5.123247

Table 4



» TODO: discuss Table 4 and add caption

» TODO: discuss Figure 9 AND better caption. May need to ask stephen for clarification on these explanations.



The graphs in Figure 9 are used to represent the accuracy of the orders of magnitude of error for using the linear interpolation method for the sample dataset. We can see that linear interpolation tends to produce errors in the range of plus or minus 2.5. There are some outliers outside of plus or minus 2.5 with worse results such as those seen in mu = 0.24. In mu = 0.24 we can see errors in the range of -5 up to 7.5 orders of magnitude. There are options to improve the errors found in linear interpolation such as sampling more points, but taking these extra samples takes the computation time of

running the CFD models that this work is looking to improve upon with the use of neural networks.

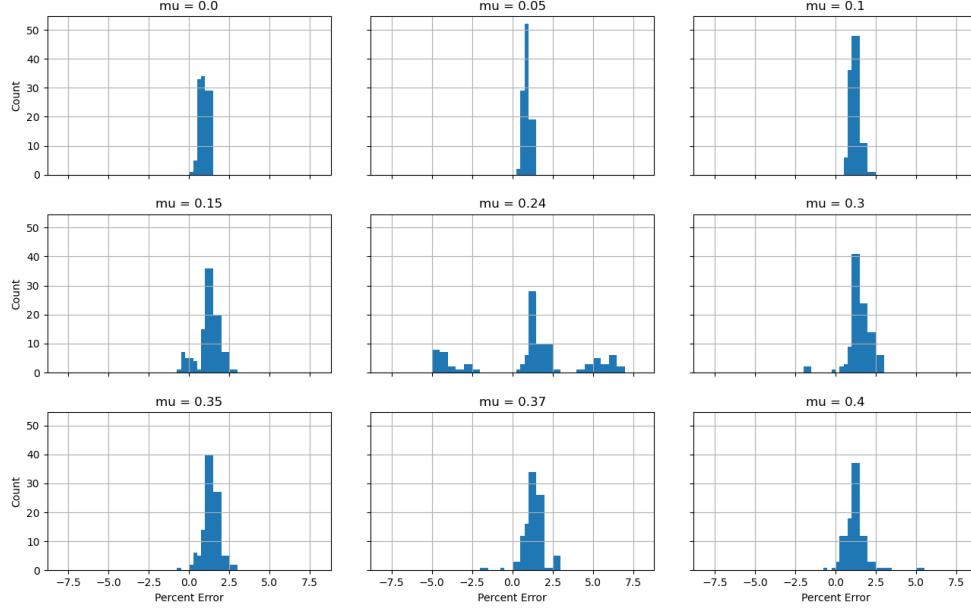


Figure 9: Percent Error Hist All interp 1

### 3.4.2 Initial Neural Network for Base Model



- » TODO: discuss Table 5 AND add caption. Again for this table, like the linear interpolation, i have no idea what it is showing. Ask stephen and if no answer then delete.
- » TODO: discuss Figure 10 AND better caption

<b>mu</b>	<b>0.0</b>	<b>0.05</b>	<b>0.1</b>	<b>0.15</b>	<b>0.24</b>	<b>0.3</b>	<b>0.35</b>	<b>0.37</b>	<b>0.4</b>
count	119	119	119	119	119	119	119	119	119
mean	0.019481	0.045969	0.094344	0.127665	0.084728	0.205783	0.155414	0.092191	0.067463
std	0.279356	0.234753	0.31285	0.517913	1.641484	0.567834	0.526382	0.634942	0.971601
min	-0.89689	-0.98214	-1.11886	-1.31219	-5.93081	-2.12056	-1.22168	-2.76316	-1.92399
25%	-0.05203	-0.06545	-0.05823	-0.10428	-0.1226	-0.08305	-0.19554	-0.22951	-0.27194
50%	0.035129	0.058321	0.071778	0.063685	0.06694	0.201968	0.133247	0.079503	0.016649
75%	0.083077	0.177541	0.219738	0.438618	0.426983	0.422871	0.399985	0.379661	0.287101
max	0.770263	0.611293	1.090276	1.656562	7.114196	2.373645	1.715405	2.491554	8.722475

Table 5

With the use of neural networks, we can see an improvement in the number of points with an order of magnitude of error found near zero. These results in Figure 10 show that the neural network approach is going to consistently provide estimates for the power

coefficient tables that will be more accurate than the values produced using the linear interpolation approach (Figure 9).

«< TODO: What is the impact on power? must ask Bob and Luke on how to word the impact. »>

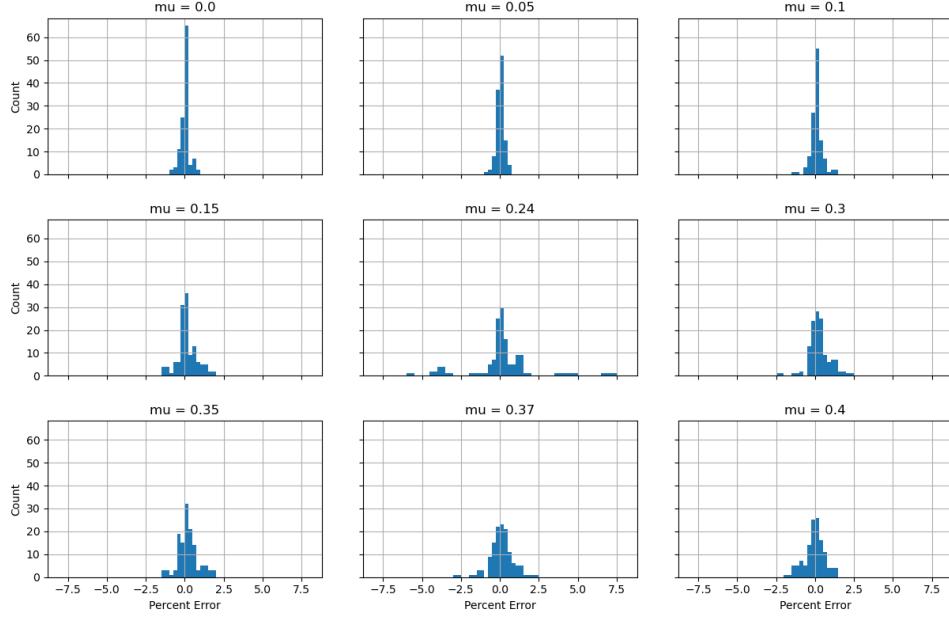


Figure 10: Percent Error Hist All samex 1

### 3.4.3 Hyperparameter Optimization Spreadsheet of Modifications

The spreadsheet in Figure 11 compiles all the results from parameter sweeps performed by Optuna, from the base model and all the intermediate models produced before arriving at the optimal RNN. Each row displays metrics including  $R^2$ , RMSE, MAPE, and the percentage of data points under a certain threshold of error. The rightmost column contains notes of the intermediate modifications to the model. The results of the final model before we tried ensembling methods are significantly better in all percent error thresholds (i.e. 20%, 10%, 5%, 2.5%, 1%, 0.5%).

ML Rotor Models (Neural Networks)											
Layer Sizes	Epochs / Optimizer			Percent Difference						Notes	
		R^2	RMSE	< 20%	< 10%	< 5%	< 2.5%	< 1%	< 0.5%		
324, 222, 126, 57, 32, 16, 9, 1 (NN5Ad)	200 / Adamax	0.99961	0.002476	0.4957%	100.0%	99.4%	99.2%	94.6%	79.0%	52.2%	> Use Adamax instead of Nadam
324, 222, 126, 57, 32, 16, 9, 1 (NN5Nad)	200 / Nadam	0.99949	0.002829	0.5290%	100.0%	99.4%	98.2%	95.6%	76.4%	53.8%	
335, 335, 139, 84, 21, 15, 6, 1	200 / Adamax	0.99956	0.002640	0.5180%	99.9%	99.6%	98.7%	95.0%	76.4%	48.0%	> 1 Unstable model training
218, 109, 98, 37, 19, 12, 5, 1 (NN4)	200 / Nadam	0.99945	0.002958	0.5928%	99.8%	99.5%	97.8%	92.6%	73.6%	47.9%	
321, 265, 132, 55, 21, 14, 6, 1	200 / Adamax	0.99958	0.002579	0.6248%	99.9%	99.6%	98.5%	94.0%	65.4%	34.8%	
											> Hidden layers: 5 -> 6 > Epochs: 0.001 -> 0.0005 > Range of input layer: [128, 400] > Use Nadam instead of Adam > Epochs determined from EarlyStopping
109, 98, 37, 19, 12, 5, 1 (NN3)	200 / Nadam	0.99944	0.002975	0.6989%	100.0%	99.8%	97.6%	91.4%	65.4%	39.4%	
109, 98, 37, 19, 12, 5, 1 (NN3)	100 / Nadam	0.99931	0.003311	0.7720%	99.8%	99.0%	96.6%	87.4%	66.6%	39.4%	
109, 90, 38, 18, 9, 5, 1 (NN2)	100 / Nadam	0.99908	0.003815	1.0092%	100.0%	99.6%	96.0%	83.6%	41.6%	21.0%	
59, 54, 26, 12, 5, 3, 1 (NN1)	100 / Nadam	0.99905	0.003867	1.0595%	99.6%	98.6%	93.2%	82.2%	47.0%	29.0%	
109, 98, 37, 19, 12, 5, 1 (NN3)	50 / Adam	0.99912	0.003727	0.9110%	100.0%	98.8%	95.8%	86.8%	52.8%	28.2%	
109, 90, 38, 18, 9, 5, 1 (NN2)	50 / Adam	0.99921	0.003542	1.0108%	99.8%	99.0%	94.0%	84.8%	51.0%	29.0%	
											> Use ReLU instead of mish > Hidden layers: 4 -> 5 > Range of input layer: [32, 128]
59, 54, 26, 12, 5, 3, 1 (NN1)	50 / Adam	0.99986	0.004254	1.3402%	98.8%	97.0%	87.2%	70.8%	42.4%	20.4%	
64, 64, 32, 11, 4, 1	50 / Adam	0.99968	0.00423	1.2034%	99.8%	99.2%	92.4%	81.2%	46.4%	24.0%	
64, 64, 24, 4, 2, 1	30 / Adam	0.99973	0.004505	1.1666%	99.8%	98.2%	89.6%	75.4%	46.8%	27.4%	
56, 56, 32, 11, 2, 1	30 / Adam	0.99968	0.004574	1.5031%	99.8%	98.8%	89.8%	63.6%	31.6%	17.2%	
51, 51, 32, 10, 4, 1	30 / Adam	0.99955	0.004790	1.4898%	99.6%	95.6%	86.8%	72.4%	34.0%	17.4%	
57, 55, 33, 13, 26, 1	30 / Adam	0.99959	0.004715	1.5823%	99.6%	98.0%	91.6%	71.2%	35.2%	16.0%	
62, 26, 13, 9, 2, 1	30 / Adam	0.99914	0.005428	1.6624%	99.0%	96.0%	86.0%	66.6%	30.6%	16.8%	
57, 53, 27, 11, 7, 1	30 / Adam	0.99964	0.004639	1.4702%	99.8%	97.4%	86.6%	67.4%	32.2%	17.2%	
50, 31, 22, 13, 2, 1	30 / Adam	0.99789	0.005774	1.4567%	99.8%	97.4%	88.4%	75.0%	39.4%	18.2%	
57, 61, 26, 7, 25, 1	30 / Adam	0.99858	0.004735	1.7875%	100.0%	96.2%	88.6%	61.4%	27.2%	13.8%	
59, 53, 23, 12, 6, 3, 1	30 / Adam	0.99828	0.005222	1.3551%	99.2%	96.8%	90.2%	74.4%	42.6%	22.4%	
65, 51, 25, 12, 5, 4, 1	30 / Adam	0.99842	0.004996	1.5995%	100.0%	98.2%	88.4%	62.2%	30.4%	15.8%	
61, 58, 25, 16, 6, 3, 1	30 / Adam	0.99876	0.004421	1.4780%	99.8%	97.6%	89.2%	73.6%	33.6%	17.2%	
											> Start Optuna HP tuning (3-6 layers) > Range of input layer: [14, 64] > Next layers between 0.5x to 1x - 1 of previous layer
62, 50, 23, 11, 5, 3, 1	30 / Adam	0.99861	0.004687	1.2966%	99.8%	97.6%	90.0%	78.4%	45.0%	23.2%	
											> Hidden layers: 3 -> 4 > Double size of input layer
32, 32, 16, 16, 2, 1	25 / Adam	0.98560	0.014880	2.0958%	99.8%	97.2%	81.4%	51.8%	23.0%	13.4%	> Use linear output instead of softmax
16, 16, 8, 4, 1 (Base Model)	25 / Adam	0.99287	0.010761	5.1823%	92.2%	68.2%	41.4%	28.4%	11.8%	6.2%	> Original model, mish activation and softmax output

Figure 11: Full Model Results Table

As mentioned, the "Notes" column in Figure 11 refers to major modifications to the network configuration discovered by Optuna. Figure 12 shows a summary of all the intermediate modifications we made to the base model to arrive at the final model. Five-hundred points were randomly sampled for clarity of visualization for each graph. We can see in the plot labeled "Base Model" in Figure 12 that the base model operated at around and above an order of magnitude of error across all angles of attack (AoA). After examining the results from the base model, the initial modification included increasing the number of hidden layers from 3 to 4, doubling the size of the input layer from 16 to 32 neurons, and using a linear activation in the output layer in place of the softmax activation. The error after the first modification shows that the model performed better with predicted values lying mostly within an order of magnitude of error (plot labeled "Modification 1" in Figure 12). The points were mostly in the positive range of an order of magnitude off of the target value. There were less points beyond an order of magnitude in error when compared to the base model. The second modification included the introduction of hyperparameter optimization with Optuna. We began tuning the number of neurons in the input layer by searching in a range of [14, 64] neurons, setting up the criteria to reduce each proceeding layer by 0.5-1 times the size of the previous layer. After the second modification, we can see that the model performed even better, with most of the points being distributed within an order of magnitude (plot labeled "Modification 2" in Figure 12). The error range for points more closely fits around no error rather than favoring the positive or negative side of an order of magnitude. We then made the decision to use the ReLU activation in the hidden layers instead of Mish. We increased the number of hidden layers from 4 to 5 and searched the range of [32, 128]

for the number of neurons in the input layer, which we labeled the third modification. After the third modification, we can see that the points lie within an order of magnitude in error, indicating that the predictions are getting better (plot labeled "Modification 3" in Figure 12). To further improve model performance, we decided to continue tuning the hyperparameters in Optuna. Based on the Optuna results, we increased the number of hidden layers from 5 to 6, decreased the learning rate from 0.001 to 0.0006, modified the number of neurons in the input layer to the range [128, 400], used the Nadam optimizer in place of Adam, and let early stopping determine the number of epochs. After this fourth modification, we see that the expected error range for the points is decreasing (plot labeled "Modification 4" in Figure 12). Finally, through hyperparameter optimization, we made the choice to use the Adamax optimizer instead of Nadam, resulting in the final model. In the final model (plot labeled "Final Model" in Figure 12), we see that model is showing the lowest expected margin of error for the prediction of the points.

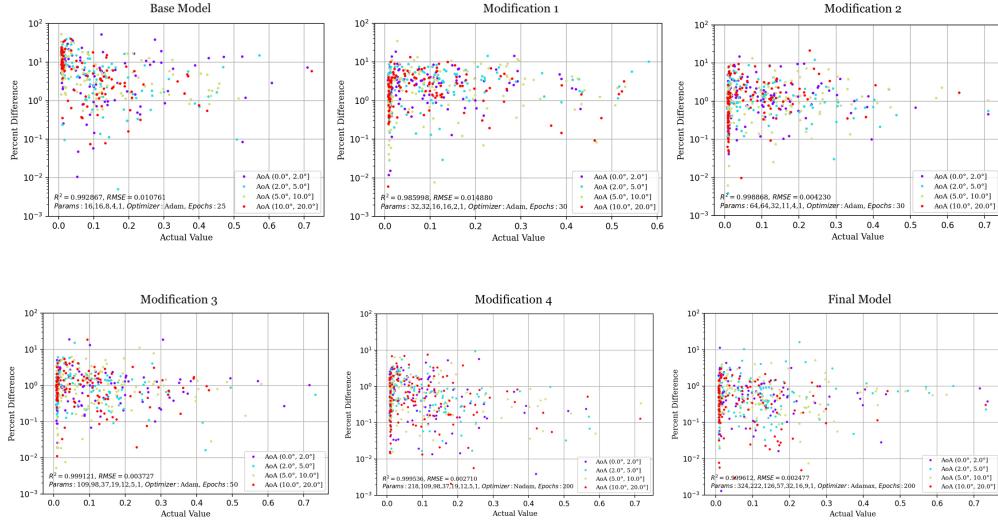


Figure 12: Full Test Results with Modifications

### 3.4.4 Final Model Before Ensemble Method

Table 6 shows the transition from the base model to the final model, specifically focusing on the choices for the hyperparameters of each model. The choices for the final model were determined through hyperparameter optimization by Optuna. The reasoning behind the adjustment of each hyperparameter from the initial choice for the base model to the final choice has been discussed in Section 3.3.

<b>Summary</b>	<b>(Base → Final)</b>
Hidden layer activation:	Mish → ReLU
Optimizer:	Adam → Adamax
Learning rate:	0.001 → 0.0006
Number of hidden layers:	3 → 6
Number of neurons per layer:	[16, 16, 8, 4, 1] → [324, 222, 126, 57, 32, 16, 9, 1]
Epochs:	25 → 200

Table 6: Model Configuration Parameters: Base Model to Final Model

Table 7 shows a distribution comparing the percent difference performance of the final model against the base model. The " $<20\%$ " column indicates the percentage of data points predicted with under a 20% error, the " $<10\%$ " column indicates the percentage of data points predicted with under a 10% error, and so on. The final model shows significant improvement over the base model in accuracy of points with a 100% likelihood that a point will never be beyond 20% error. The final model also has the ability to predict points within 5% error 99.2% of the time which is a significant improvement over the base model only being accurate within 5% error 41.4% of the time. Table 7 shows that the implementation of hyperparameter optimization improved model performance.

<b>Model</b>	<b><math>&lt;20\%</math></b>	<b><math>&lt;10\%</math></b>	<b><math>&lt;5\%</math></b>	<b><math>&lt;2.5\%</math></b>	<b><math>&lt;1\%</math></b>	<b><math>&lt;0.5\%</math></b>
Base	92.20%	68.20%	41.40%	28.40%	11.80%	6.20%
Final	100.00%	99.40%	99.20%	94.60%	79.00%	52.20%

Table 7: Percent Difference Thresholds (Base vs. Final)

### 3.5 Ensemble Methods

After the initial success of introducing a neural network into our problem and implementing hyperparameter optimization, we decided to explore additional methods to further improve model performance of the final model discussed in the previous section. We considered ensemble methods. Ensembling is an optimization method often used in top-performing models for machine learning problems [7]. The idea behind ensembling is to combine the predictions of multiple models to leverage the differing strengths of each model for more robust predictions. Ensembling can take the form of a simple statistical analysis, such as taking the average or median of predictions. Alternatively, it can involve stacking, where multiple models are concatenated to form a new model that is designed to learn how to best combine the stacked predictions. There are several different kinds of ensembling methods, and we are going to briefly introduce a few of them in this section.

Concatenation is the process of combining models side-by-side, which can suffer from an explosion of dimensionality as a result. The simple average (or mean) method is performed by taking the average predictions of each model. Unlike concatenation, simple average does not suffer from increased dimensionality problems; it makes the assumption that the data is normally-distributed. The weighted average method, which is another

option, takes tensor outputs and multiplies the outputs by weights, then linearly combines the weights. The total of the weights multiplied by all models must add to one for the weighted average method to work. The weighted average method allows the network's designer to designate the contribution of each model to the final prediction by altering the weights.

In the following sections we are going to discuss the results of applying neural network ensembling methods (Section 3.5.1), such as the simple average and weighted average, using the best models from Optuna. We will follow the Optuna results with a look at gradient boosting ensembling methods (Section 3.5.2) using a well known library called XGBoost.

### 3.5.1 Neural Network Ensembles

**3.5.1.1 Simple Average** Figure 13 is the plot of the best neural network found by Optuna before we considered any ensemble methods. The simplest ensemble method is to take the average of multiple neural network predictions. Figure 14 shows the models chosen for the ensemble. The three chosen models are within the top five models found in Optuna.

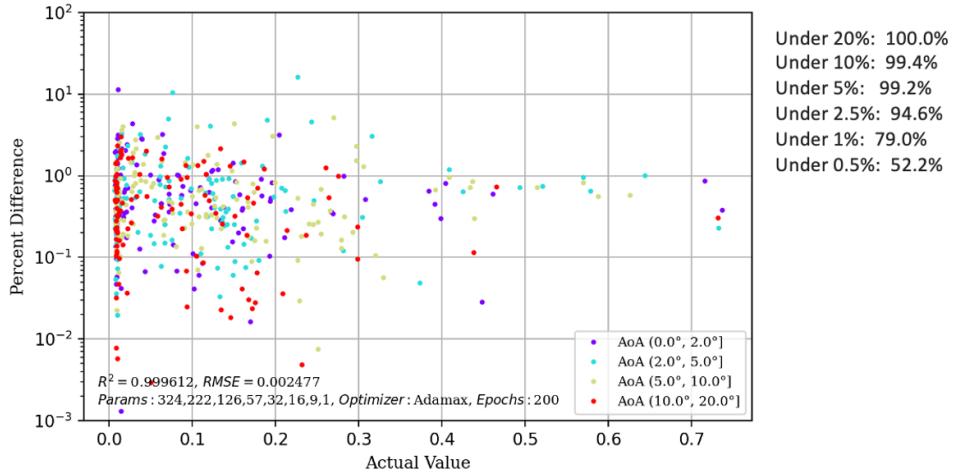


Figure 13: Sample output from best neural network before ensembling

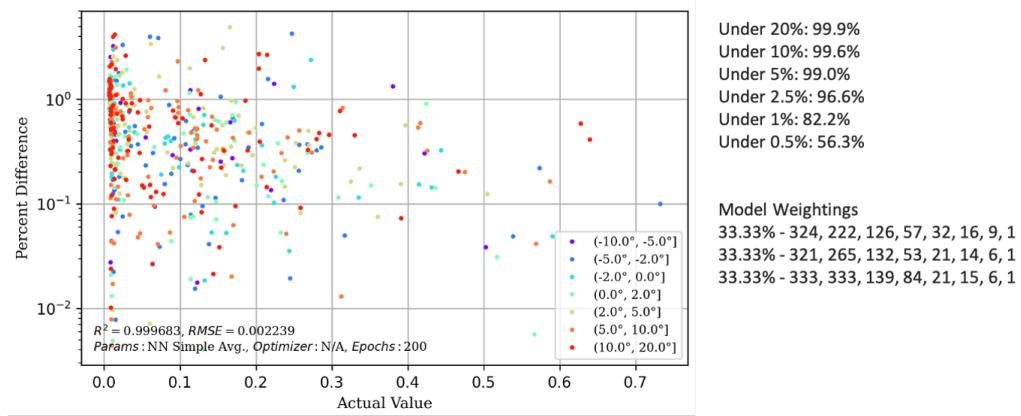


Figure 14: Neural Network with Simple Average Ensemble Method

In Figure 14 we already see a noticeable improvement over the single neural network model in Figure 13, with higher  $R^2$ , lower RMSE, and a higher percentage of points predicted under 2.5%, 1%, and 0.5% error.

**3.5.1.2 Weighted Average** A weighted average allows us to place more or less weight on the predictions of certain neural networks. A Cartesian product grid search was conducted to test possible weight vectors, and the optimal combination ended up being 31.25%, 56.25%, 12.5%. However, compared to the simple average plot and error percentages in Figure 14, there was almost no improvement using the weighted average (Figure 15).

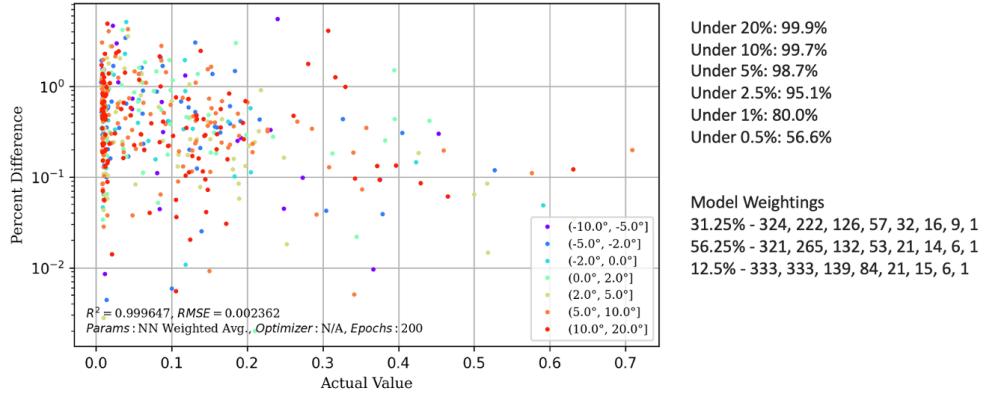


Figure 15: Neural Network with Weighted Average Ensemble Method

### 3.5.2 Gradient Boosting Ensembles

While applying ensembling methods to the network models created by Optuna, Gradient Boosting Ensembling also came into our discussion. Gradient boosting is a more classical machine learning technique that has been known to work well for most problems [6, 13]. It is a type of ensembling as each iteration a new decision tree learns to correct mistakes from previous decision trees [1]. Its advantage over a neural network is that it has a simpler structure and is much faster to train, which makes this method a choice when the dataset is not large (as in our case). The most popular gradient boosting libraries are XGBoost, CatBoost, and LightGBM. For this work, we ended up using XGBoost based on comparison of results using the Python library known as AutoML [15]. Even though XGBoost is the slowest of the three gradient boosting libraries mentioned here, the difference in training speed is negligible for our problem.

**3.5.2.1 XGBoost Initial Model** The default XGBoost model was initially run on our dataset with the parameters defined in the table below. Max Depth is the maximum depth of a tree. Increasing this value makes the model more complex and more prone to overfitting. N\_Estimators is the number of gradient boosted trees. Learning Rate is the boosting learning rate. Max Bin is the maximum number of discrete bins to bucket continuous features. Subsample is the subsample ratio of the training instances. Setting it to 0.5 means that XGBoost would randomly sample half of the training data prior to growing trees. Lambda is the L2 regularization term on weights [8].

Max Depth	6
N_Estimators	100
Learning Rate	0.3
Max Bin	256
Subsample	1
Lambda	0

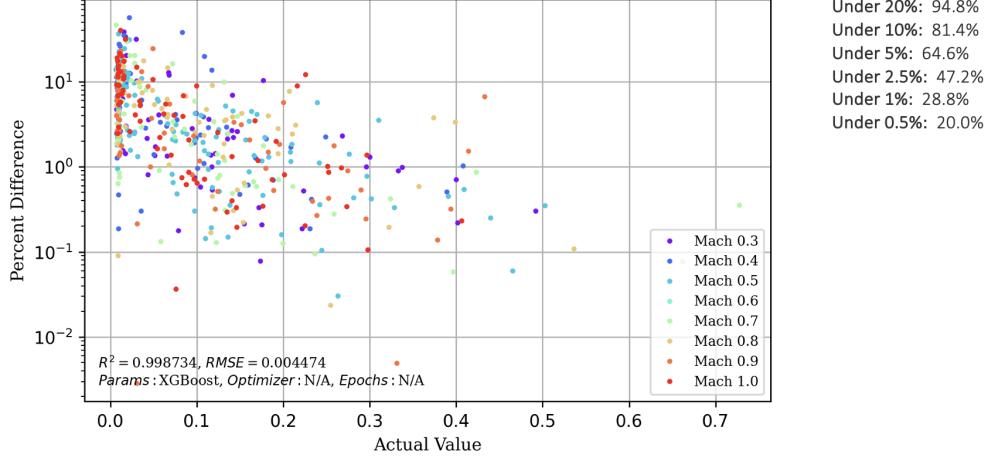


Figure 16: XGBoost Base Model

The base model only took a few seconds to train, but we observe in Figure 16 that the performance is unsatisfactory, as there are a large cluster of points above 10% error. We needed to tune the model to see how promising it was for locating an improvement over the best neural network model before applying ensembling.

**3.5.2.2 Single Best XGBoost** Optuna is used for automating the optimization process of these hyperparameters modified by XGBoost. These initial conditions were specified for Optuna to conduct a hyperparameter search:

Max Depth	5 to 20
Learning Rate	0.01 to 0.3, step 0.01
N_Estimators	1000, with early_stopping_rounds of 10 iterations
Subsample	0.4 to 1.0, step 0.05
Alpha, Lambda, Gamma	0 to 5

Alpha is the L1 regularization term on weights, and gamma is the minimum loss reduction required to make a further partition on a leaf node of the tree. The larger gamma is, the more conservative the algorithm will be [8]. From the slice plot of the optimization trials, all the best models had an alpha and gamma value of 0, so these two parameters were eliminated from the search. In addition, all the top models had a subsample rate of at least 0.6, so the range was reduced accordingly:

Max Depth	5 to 20
Learning Rate	0.01 to 0.3, step 0.01
N_Estimators	1000, with early_stopping_rounds of 10 iterations
Subsample	0.6 to 1.0, step 0.05
Lambda	0 to 5

After tuning our model with Optuna, we see an order of magnitude improvement in the results in Figure 17 compared to the XGBoost base model results in Figure 16. Most points are now under 10% error, and its error ranges are on-par with some deeper neural networks. Its  $<20\%$ ,  $<10\%$ , and  $<5\%$  error percentages are worse than several neural networks. However, its  $<1\%$  error percentage is similar to the 4th best NN, and the  $<0.5\%$  error percentage is close to the 3rd best NN. While considerably more complex than the base model, this tuned model took less than two minutes to train.

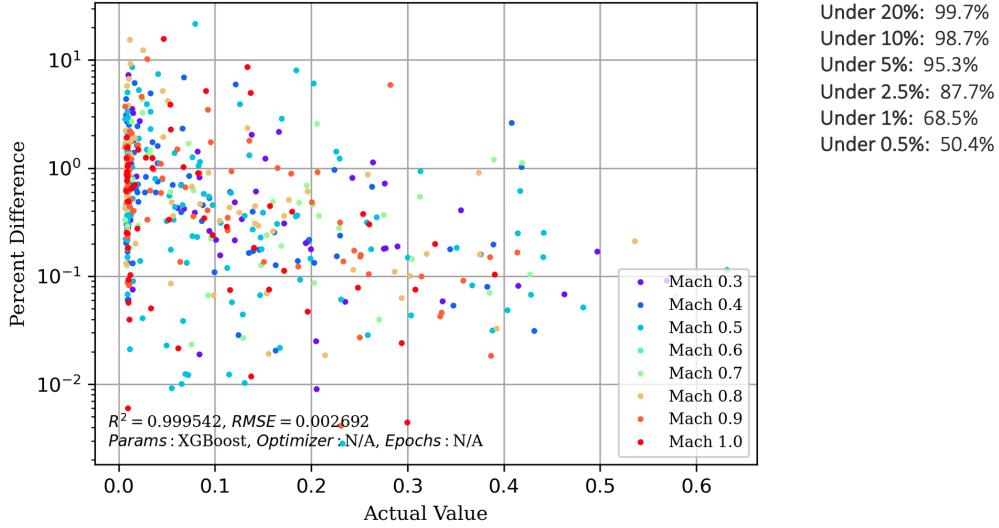


Figure 17: Single Best XGBoost Model Results

**3.5.2.3 Voting Regressor (Top 3)** The simplest ensemble model for XGBoost is using’s scikit-learn’s VotingRegressor, which takes several base estimators and then simply averages the individual predictions to form a final prediction. For the base regressors, we simply took the top three XGBoost models from Optuna’s hyperparameter tuning. Compared to the best individual XGBoost model in Figure 17, we observe another noticeable improvement in the results using VotingRegressor (Figure 18), with higher percentages for the error ranges  $<5\%$ ,  $<2.5\%$ ,  $<1\%$ , and  $<0.5\%$ .

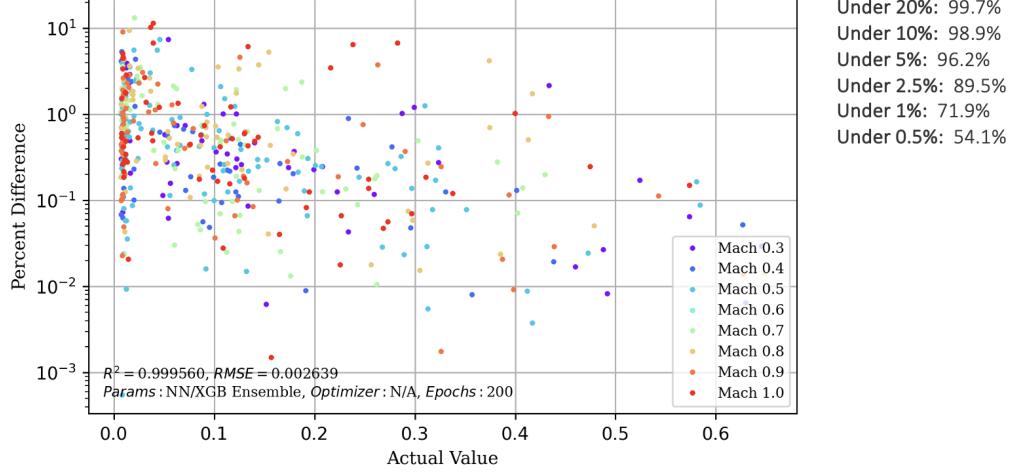


Figure 18: VotingRegressor Results

**3.5.2.4 Weighted Average** VotingRegressor provides an additional parameter called 'weights' should we want to take a weighted average of the XGBoost predictions. However, the best weights found from grid search produced no improvement. The plot is not shown here, but the results from this test are included in Figure 22.

### 3.5.3 Hybrid Ensembles

It is recommended that the VotingRegressor take conceptually different machine learning models as base regressors. In Figure 18 we were using three XGBoost models, which are conceptually very similar. Scikit-learn does not provide an easy way to combine neural networks and XGBoost. However, as a baseline, we can do similar statistical analysis like we did for neural networks alone. Figures 19 and 20 show the results for the average and median of predictions from the top 5 neural networks (Table 8) and top 3 XGBoost models (Table 9) from Optuna, resulting in a hybrid ensemble model:

Parameters	Optimizer
324, 222, 126, 57, 32, 16, 9, 1	Adamax
324, 222, 126, 57, 32, 16, 9, 1	Nadam
333, 333, 139, 84, 21, 15, 6, 1	Adamax
218, 109, 98, 37, 19, 12, 5, 1	Nadam
321, 265, 132, 53, 21, 14, 6, 1	Adamax

Table 8: Top 5 Neural Network Model Configurations (LR = 0.0006, 200 epochs each)

Max Depth	Learning Rate	Bagging Fraction	Random State
20	0.0138	0.867	2020
19	0.017	0.867	48
20	0.015	0.86	202

Table 9: Top 3 XGBoost Model Configurations ( $n_{\text{estimators}} = 750$ )

While the average has its uses, taking the median may be more robust, since it is not highly affected by outliers from skewing results. In the rare event that one of the models fails to train, the result will also not be ruined. In both cases, the ensemble outperforms both neural network and XGBoost-only ensembles, as we see another large improvement in percentages for  $< 2.5\%$ ,  $< 1\%$ , and  $< 0.5\%$  error. The median (Figure 20) does slightly better than a simple average (Figure 19). So, the overall best model is the median hybrid ensemble.

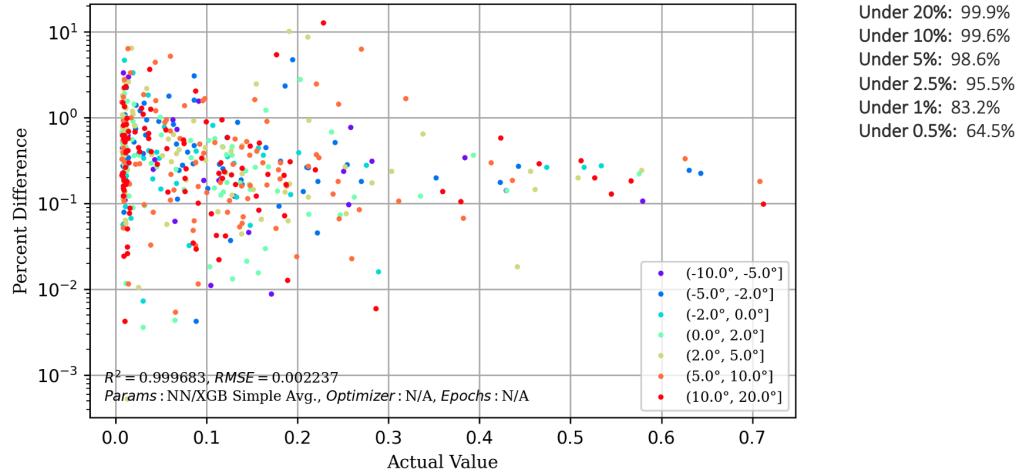


Figure 19: Neural Network + XGBoost Hybrid Average Ensemble Results

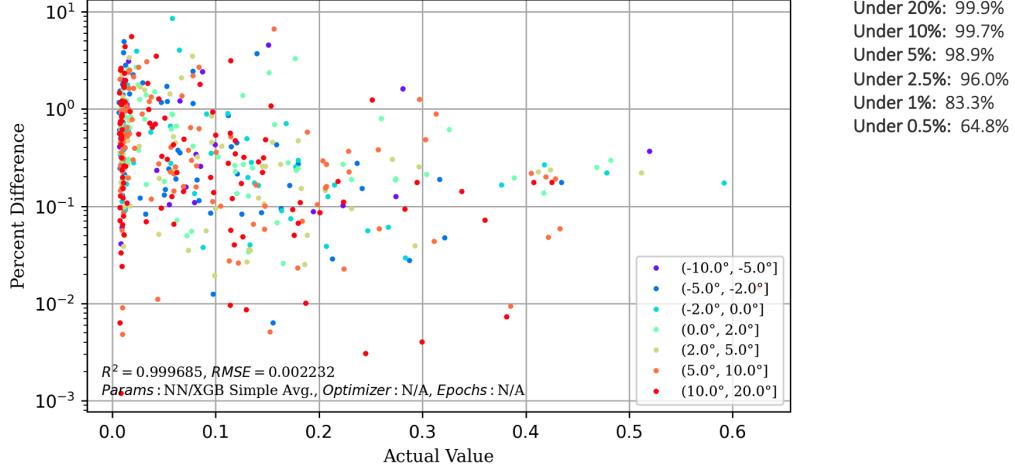


Figure 20: Neural Network + XGBoost Hybrid Median Ensemble Results

### 3.6 Summary of Optuna Hyperparameter Optimization, XGBoost, and Ensembling Methods

We close the results section with a description and summary tables showing the best results from neural network modeling, XGBoost, and ensembling methods. The first table we are going to look at is the neural network results (Figure 21). We can see a high degree of improvement as we moved from the base model, transitioning to models created by hyperparameter optimization, and lastly to ensemble models created by applying simple average, weighted average, and median ensembling methods leveraging the best models located by Optuna. The results of this table show that as we move from the base hand-generated model to applying hyperparameter optimization and then to ensembling methods, we see significant improvements in the number of points within each error range. A notable observation is that the best model created by ensembling with neural networks was able to generate 95.7% of the points within 2.5% error.

ML Rotor Models (Neural Networks)	Epochs / Optimizer	R^2	RMSE	MAPE	Percent Difference					Notes
					< 20%	< 10%	< 5%	< 2.5%	< 1%	
NN Ensemble (Weighted Avg. 2)	200 / Adamax + Nadam	0.99957	0.002595	N/A	99.9%	99.7%	98.8%	95.7%	82.7%	> Same models as Simple Avg. 2
NN Ensemble (Simple Avg. 2)	200 / Adamax + Nadam	0.99957	0.002598	N/A	99.9%	99.7%	98.8%	95.5%	82.7%	> Weighting: 0.25, 0.375, 0.375 > Optimizer: Adamax + Nadam > Weighting: 0.25, 0.375, 0.375 > Optimizer: Adamax + Nadam
NN Ensemble (Weighted Avg.)	200 / Adamax	0.99965	0.002362	N/A	99.9%	99.7%	98.7%	95.1%	80.0%	> Same models as Simple Avg. 1
NN Ensemble (Median)	200 / Adamax	0.99966	0.002304	N/A	99.9%	99.6%	98.9%	96.2%	80.8%	> Weighting: 0.3125, 0.5625, 0.125
NN Ensemble (Simple Avg.)	200 / Adamax	0.99968	0.002239	N/A	99.9%	99.6%	99.0%	96.6%	82.2%	> Optimizer: Adamax instead of Nadam
324, 222, 126, 57, 32, 16, 9, 1 (NN5Ad)	200 / Adamax	0.99961	0.002476	100.0%	99.4%	99.2%	94.6%	79.0%	52.2%	> Use Adamax instead of Nadam
324, 222, 126, 57, 32, 16, 9, 1 (NN5Nad)	200 / Nadam	0.99949	0.002829	0.5290%	100.0%	99.4%	98.2%	95.6%	76.4%	53.8%
333, 333, 139, 84, 21, 15, 6, 1	200 / Adamax	0.99956	0.002640	0.5180%	99.9%	99.6%	98.7%	95.0%	76.4%	48.0% > Unstable model training
218, 109, 98, 37, 19, 12, 5, 1 (NN4)	200 / Nadam	0.99945	0.002958	0.5928%	99.8%	99.5%	97.8%	92.6%	73.6%	47.9%
321, 265, 132, 53, 21, 14, 6, 1	200 / Adamax	0.99958	0.002579	0.6248%	99.9%	99.6%	98.5%	94.0%	65.4%	34.8%
16, 16, 8, 4, 1 (Base Model)	25 / Adam	0.99287	0.010761	5.1823%	92.2%	68.2%	41.4%	28.4%	11.8%	6.2% > Original model, mish activation and softmax output

Figure 21: Neural Network Summary of Results

The results of working with XGBoost is seen in Figure 22 found below. The initial

XGBoost model performance is viewable on the last line of the figure. After the initial model was created, the team followed the same approach we had done when designing the neural networks by applying Optuna to make architectural configuration choices. Optuna produced several XGBoost models with increasingly accurate results. We can observe that Optuna started its approach by increasing the depth from 6 to 13 and eventually settled around 19 to 20 as the ideal depth for the XGBoost network. We can also see an immediate improvement over the base XGBoost model to the models using Optuna with a 5% increase in the number of points within 20% error. There was also a significant increase in the number of points within 10%, 5%, 2.5%, 1%, and 0.5% error. Results of the points within the targeted error ranges continued to marginally increase throughout each Optuna iteration featured in the table. It should be noted that XGBoost did not produce a better model every iteration, but tried many different combinations of hyperparameters to find the networks which identified points within smaller regions of error. After applying Optuna to the XGBoost model, we had the result of 96.2% of the points within 5% error, but we were still looking to improve the model. We then applied the VotingRegressor and weighted average ensemble methods (described in Sections 3.5.2.3 and 3.5.2.4 respectively) using the best models produced from XGBoost with Optuna runs.

ML Rotor Models (Gradient-Boosting) Model	max_depth	n_estimators	lr	max_bins	bag_fractions	lambda	R^2	RMSE	Percent Difference					Notes
									< 20%	< 10%	< 5%	< 2.5%	< 1%	
Ensemble (Weighted Avg.)	N/A	N/A	N/A	N/A	N/A	N/A	0.99956	0.002639	99.7%	98.9%	96.2%	89.5%	72.0%	54.1% > XGB 1, 2, 3: (0.28, 0.36 0.36)
Ensemble (Voting Reg.)	N/A	N/A	N/A	N/A	N/A	N/A	0.99956	0.002639	99.7%	98.9%	96.2%	89.5%	71.9%	54.1% > Voting regressor using XGB 1, 2, 3
XGB Opt 1	20	750	0.0138	372	0.867	0	0.99954	0.002692	99.7%	98.7%	95.3%	87.7%	68.5%	50.4% > With Optuna HP Tuning
XGB Opt 2	19	750	0.0175	277	0.867	0	0.99954	0.002689	99.7%	98.7%	95.5%	87.8%	69.0%	50.9%
XGB Opt 3	20	750	0.0153	275	0.860	0	0.99954	0.002694	99.7%	98.7%	95.5%	87.7%	68.9%	50.4%
XGB Base Model	6	100	0.3	256	1	0	0.99873	0.004484	94.8%	81.4%	64.6%	47.2%	28.8%	20.0%

Figure 22: Gradient-Boosting Summary of Results

The final figure discussed here is Figure 23 which shows what we are calling the hybrid ensemble model. The hybrid ensemble model is a combination of the best neural networks and XGBoost models prior to ensembling. It was interesting to combine the benefits of neural networks and decision tree support to get a better result than either individual method was able to accomplish.

Layer Sizes	Epochs / Optimizer	R^2	RMSE	MAPE	Percent Difference					Notes
					< 20%	< 10%	< 5%	< 2.5%	< 1%	
NN/XGB 8-Ensemble (Median)	200 / Adamax + Nadam	0.99968	0.002231	N/A	99.9%	99.7%	98.9%	96.0%	83.3%	64.8% > XGB: Opt 1, 2, 3
NN/XGB 8-Ensemble (Simple Avg.)	200 / Adamax + Nadam	0.99968	0.002237	N/A	99.9%	99.6%	98.6%	95.5%	83.2%	64.5%
Initial NN/XGB 8-Ensemble (Median)	50 / Adamax + Nadam	0.99963	0.002433	N/A	99.9%	99.5%	98.1%	94.0%	77.6%	58.1% > NN: 5Adamax, 5Nadam, 6, 7Adamax, 7Nadam

Figure 23: Hybrid Ensembles Summary of Results

## 4 Conclusion

Due to the computational time requirements of running CFD models for estimating the rotorcraft airfoil performance being restrictive, having a ML surrogate model would allow for estimating a large number of potential airfoil designs in a much shorter time-frame. In this work we provided an analysis of the effects of applying hyperparameter optimization

to the base ML surrogate model and employing ensembling to improve accuracy, seeking an alternative for estimating the airfoil tables that were previously produced by running CFD simulations. The motivation behind this effort was based on the computational expense of using CFD models to generate the airfoil performance tables. We began by utilizing linear interpolation between CFD points to estimate the values for an airfoil table, using these values to approximate the power coefficient of the CFD model for modification to the rotorcraft airfoil geometry. Though linear interpolation proved to be computationally efficient and reasonably accurate, the benefits of this approach begin to decrease on large, irregularly sampled parameter spaces. A RNN was manually designed to combat the accuracy issues of linear interpolation, but the base model proved to be insufficient on the provided dataset and suffered from overfitting. Upon investigating approaches to improve on the base model, we decided to explore hyperparameter optimization using the Optuna library. The addition of hyperparameter tuning to find the optimal values for our neural network architecture resulted in a significant improvement in model performance over the base model. We next began exploring ensemble methods to build on the success of hyperparameter optimization, leading to an analysis of neural network and gradient boosting ensembling methods. We performed tests using the simple average and weighted average ensembling schemes for our top-performing neural networks. Though both schemes provided similar results, we did observe a noticeable improvement in model performance over our Optuna-designed neural network. Continuing our exploration of ensemble methods we began utilizing XGBoost for ensembling using decision trees. We manually designed an XGBoost model that yielded unsatisfactory results, so we modified the hyperparameters in our XGBoost model using Optuna, which improved performance. We used the VotingRegressor and weighted average ensembling schemes on our top 3 Optuna-designed XGBoost models, both tests producing similar results but still an improvement over the best XGBoost model chosen by Optuna. Finally, we combined the top 5 neural networks and top 3 XGBoost models designed using Optuna into a hybrid ensemble model, using the simple average and median ensembling schemes. The resulting models from both of these tests proved to perform better than all of our other models, with the median hybrid ensemble model being determined to be the best performing model due to its robustness.

The results of this study yielded a ML model for estimating the performance of an aircraft airfoil that was able to decrease analysis time and improve accuracy of predicting air performances tables. Future work for this effort includes applying this workflow to the lift and moment .....

## References

- [1] Aliyev, Vagif. *Gradient Boosting Classification explained through Python*. Oct. 2020. URL: <https://towardsdatascience.com/gradient-boosting-classification-explained-through-python-60cc980eeb3d> (visited on 09/13/2022).
- [2] Allen, Luke D. et al. “Rotor Blade Design Framework for Airfoil Shape Optimization with Performance Considerations”. In: *AIAA Scitech 2021 Forum*. DOI: [10.2514/6.2021-0068](https://arc.aiaa.org/doi/abs/10.2514/6.2021-0068). URL: <https://arc.aiaa.org/doi/abs/10.2514/6.2021-0068>.
- [3] Amsallem, David. *Interpolation on manifolds of CFD-based fluid and finite element-based structural reduced-order models for on-line aeroelastic predictions*. Stanford University, 2010.
- [4] Brahme, Anders. *Comprehensive biomedical physics*. Newnes, 2014.
- [5] Brunton, Steven L, Noack, Bernd R, and Koumoutsakos, Petros. “Machine learning for fluid mechanics”. In: *Annual review of fluid mechanics* 52 (2020), pp. 477–508.
- [6] Chen, Tianqi and Guestrin, Carlos. “Xgboost: A scalable tree boosting system”. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 2016, pp. 785–794.
- [7] Demir, Necati. “Ensemble methods: Elegant techniques to produce improved machine learning results”. In: *Toptal Engineering Blog* (2016).
- [8] Developers, XGBoost. *XGBoost Parameters*. URL: <https://xgboost.readthedocs.io/en/stable/parameter.html> (visited on 09/13/2022).
- [9] Dogo, E. M. et al. “A Comparative Analysis of Gradient Descent-Based Optimization Algorithms on Convolutional Neural Networks”. In: *2018 International Conference on Computational Techniques, Electronics and Mechanical Systems (CTEMS)*. Dec. 2018, pp. 92–99. DOI: [10.1109/CTEMS.2018.8769211](https://doi.org/10.1109/CTEMS.2018.8769211).
- [10] Elsayed, Khairy and Lacor, Chris. “Modeling, analysis and optimization of air-cyclones using artificial neural network, response surface methodology and CFD simulation approaches”. In: *Powder technology* 212.1 (2011), pp. 115–133.
- [11] Gatski, Thomas B and Speziale, Charles G. “On explicit algebraic stress models for complex turbulent flows”. In: *Journal of fluid Mechanics* 254 (1993), pp. 59–78.
- [12] Hammond, James et al. “Machine Learning Methods in CFD for Turbomachinery: A Review”. In: *International Journal of Turbomachinery, Propulsion and Power* 7.2 (2022), p. 16.
- [13] Hancock, John T and Khoshgoftaar, Taghi M. “CatBoost for big data: an interdisciplinary review”. In: *Journal of big data* 7.1 (2020), pp. 1–45.
- [14] Hardesty, Larry. “Explained: neural networks”. In: *MIT News* 14 (2017).
- [15] He, Xin, Zhao, Kaiyong, and Chu, Xiaowen. “AutoML: A survey of the state-of-the-art”. In: *Knowledge-Based Systems* 212 (2021), p. 106622.
- [16] Jain, R. et al. “An Assessment of RCAS Performance Prediction for Conventional and Advanced Rotor Configurations”. In: *Journal of the American Helicopter Society* 61.4 (Oct. 2016), pp. 1–12. DOI: [10.4050/JAHS.61.042005](https://doi.org/10.4050/JAHS.61.042005).

- [17] Johnson, Wayne. "Rotorcraft Aerodynamics Models for a Comprehensive Analysis". In: *American Helicopter Society Annual Forum*. Washington, D.C., 1998.
- [18] Kandel, Ibrahem, Castelli, Mauro, and Popović, Aleš. "Comparative Study of First Order Optimizers for Image Classification Using Convolutional Neural Networks on Histopathology Images". In: *Journal of Imaging* 6.9 (2020). ISSN: 2313-433X. URL: <https://www.mdpi.com/2313-433X/6/9/92>.
- [19] *Keras: the Python deep learning API*. URL: <https://keras.io/> (visited on 09/06/2022).
- [20] Kingma, Diederik P. and Ba, Jimmy. *Adam: A Method for Stochastic Optimization*. 2014. URL: <https://arxiv.org/abs/1412.6980>.
- [21] Kochkov, Dmitrii et al. "Machine learning–accelerated computational fluid dynamics". In: *Proceedings of the National Academy of Sciences* 118.21 (2021), e2101784118. URL: <https://www.pnas.org/doi/abs/10.1073/pnas.2101784118>.
- [22] Koning, W. J. F., Johnson, W., and Allan, B. G. "Generation of Mars Helicopter Rotor Model for Comprehensive Analyses". In: *AHS International Technical Meeting on Aeromechanics Design for Transformative Vertical Flight*. San Fransisco, CA, Jan. 2018.
- [23] Mahgerefteh, H., Atti, O., and Denton, G. "An interpolation technique for rapid CFD simulation of turbulent two-phase flows". In: *Process Safety and Environmental Protection* 85.1 (2007), pp. 45–50.
- [24] Martins, David. *XGBoost: A Complete Guide to Fine-Tune and Optimize your Model*. Dec. 2021. URL: <https://towardsdatascience.com/xgboost-fine-tune-and-optimize-your-model-23d996fab663> (visited on 09/16/2022).
- [25] Misra, Diganta. *Mish: A Self Regularized Non-Monotonic Activation Function*. 2019. URL: <https://arxiv.org/abs/1908.08681>.
- [26] *Optuna: A hyperparameter optimization framework — Optuna 3.0.0 documentation*. URL: <https://optuna.readthedocs.io/en/stable/index.html> (visited on 09/06/2022).
- [27] Palmer, Grant. "Construction of CFD solutions using interpolation rather than computation with the ADSI code". In: *47th AIAA aerospace sciences meeting including the new horizons forum and aerospace exposition*. 2009, p. 141.
- [28] Rajagopalan, R. G. et al. "RotCFD – A Tool for Aerodynamic Interference of Rotors: Validation and Capabilities". In: *American Helicopter Society Future Vertical Lift Aircraft Design Conference*. San Fransisco, CA, Jan. 2012.
- [29] Recht, Benjamin. "A tour of reinforcement learning: The view from continuous control". In: *Annual Review of Control, Robotics, and Autonomous Systems* 2 (2019), pp. 253–279.
- [30] Saberi, H. et al. "Overview of RCAS Capabilities, Validations, and Rotorcraft Applications". In: *American Helicopter Society 71st Annual Forum*. Virginia Beach, VA, May 2015.

- [31] Srivastava, Nitish et al. “Dropout: A Simple Way to Prevent Neural Networks from Overfitting”. In: *Journal of Machine Learning Research* 15.56 (2014), pp. 1929–1958. URL: <http://jmlr.org/papers/v15/srivastava14a.html>.
- [32] Team, Keras. *Keras documentation: KerasTuner API*. URL: [https://keras.io/api/keras\\_tuner/](https://keras.io/api/keras_tuner/) (visited on 09/06/2022).
- [33] *TensorBoard*. URL: <https://www.tensorflow.org/tensorboard> (visited on 09/06/2022).
- [34] Uzair, Muhammad and Jamil, Noreen. “Effects of Hidden Layers on the Efficiency of Neural networks”. In: *2020 IEEE 23rd International Multitopic Conference (INMIC)*. 2020, pp. 1–6. DOI: [10.1109/INMIC50486.2020.9318195](https://doi.org/10.1109/INMIC50486.2020.9318195).
- [35] Vinuesa, Ricardo and Brunton, Steven L. “Enhancing computational fluid dynamics with machine learning”. In: *Nature Computational Science* 2.6 (2022), pp. 358–366.
- [36] Vinuesa, Ricardo et al. “The role of artificial intelligence in achieving the Sustainable Development Goals”. In: *Nature communications* 11.1 (2020), pp. 1–10.
- [37] Wang, Bo and Wang, Jingtao. “Application of Artificial Intelligence in Computational Fluid Dynamics”. In: *Industrial & Engineering Chemistry Research* 60.7 (2021), pp. 2772–2790. URL: <https://doi.org/10.1021/acs.iecr.0c05045>.
- [38] Wang, Bo and Wang, Jingtao. “Application of artificial intelligence in computational fluid dynamics”. In: *Industrial & Engineering Chemistry Research* 60.7 (2021), pp. 2772–2790.
- [39] Wang, Lei et al. “Radial basis function neural networks-based modeling of the membrane separation process: hydrogen recovery from refinery gases”. In: *Journal of Natural Gas Chemistry* 15.3 (2006), pp. 230–234.
- [40] Zulkifli, Hafidz. *Understanding Learning Rates and How It Improves Performance in Deep Learning*. Jan. 2018. URL: <https://towardsdatascience.com/understanding-learning-rates-and-how-it-improves-performance-in-deep-learning-d0d4059c1c10> (visited on 10/05/2022).

Placeholder for our ideas on future work. Write up abstracts for these ideas

1. global, categorical, ensemble for  $<.7$  mu

The screenshot shows a Microsoft Word document with a table titled "Mach Number". The table has "Angle of Attack" listed in the rows and Mach numbers (0.1, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1) listed in the columns. The data cells contain numerical values. The table is styled with alternating row colors (light gray and white). The Microsoft Word ribbon is visible at the top.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
	Mach Number														
1															
2															
3															
4															
5															
6															
7															
8															
9															
10															
11															
12															
13															
14															
15															
16															
17															
18															
19															
20															
21															
22															
23															
24															
25															
26															
27															
28															
29															
30															
31															
32															
33															
34															
35															
36															
37															
38															
39															
40															

Figure 24: AreaOfFocusMachNumber

2. Focus the categorial to be or minus .5 to more closely look at the problem area

## Appendix: TODO: to be deleted

Ensemble Methods Created by Ross, James E ERDC-RDE-ITL-MS CIV, last modified by Wang, Stephen ERDC-RDE-CRREL-NH CIV on May 26, 2022, viewed 31 times Types of ensemble methods: (<https://towardsdatascience.com/neural-networks-ensemble-33f33bea7df3>)

Concatenation: Combining different models side-by-side. Can cause an explosion of dimensionality Simple Average: Average the predictions of each model. No increase in dimensionality, but the average operation imposes the assumption that the data is normally-distributed. Weighted Average: Tensor outputs are multiplied by a weight and the linearly combined. Allows certain models to contribute more/less heavily to the final result. Weights must add to 1 Median: Take the median prediction for each data point. No increase in dimensionality, and more robust to outliers. Notebook with examples of types of methods: [https://github.com/cerlymarco/MEDIUM\\_NoteBook/blob/master/NeuralNet\\_Ensemble/NeuralNet\\_Ensemble.ipynb](https://github.com/cerlymarco/MEDIUM_NoteBook/blob/master/NeuralNet_Ensemble/NeuralNet_Ensemble.ipynb)

Variations on the three major themes of ensemble method

“Committee of networks” - training a collection of models with the same configuration by different initial weights Training Data: Vary the choice of data used to train each model in the ensemble. Combinations: Vary the choice of the way that outcomes from ensemble members are combined. Ensemble Models: Vary the choice of the models used in the ensemble.

«< Possible close Neural network components and selections section »> In this section we discussed the database, neural network libraries, hyperparameter optimization libraries, and the neural network component options. We also provided some insights into how we made the choices relating to these components for our problem. There are various options we could have chosen for the different components, and we will continuously revisit these choices for further optimization of our approach. In the following section, we will examine at the results of the choices we made during our examination of the components of our neural network.

## FUTURE WORK

1. Could do a write up comparing the application of the work from this paper to lift (thrust coefficient more sensitive) and moment
2. AoA range write up with added experiments based on Bob's suggestions
3. possibly use softmax for regression problem: create bins for probability distribution related to prediction error

TODO: The end result that needs to be discussed in the conclusion is from linear interpolation to optuna to xgboost to ensemble to the hybrid ensemble. Did any of our attempts beat linear interpolation? Was the hybrid ensemble the best or was optuna or xgboost by itself better? »»>

1. The results defined in this work are an analysis of the effects of modifying the hyperparameters of the base ML surrogate model seeking results comparable to the approach of CFD combined with linear interpolation.
2. In this work we started with using interpolation between CFD points in order to

estimate the values for an air foil table. These values would be used to approximate the power coefficient of the CFD model for modification to the rotorcraft airfoil geometry. briefly described issues.

3. we moved to a hand designed RNN and had these issues
4. We next moved to using hyper parameter optimization library which gave us the best results of XYZ.
5. We next moved to emsembling methods and had these results
6. Last part of ensemble was trying the hybrid ensemble
7. After all tests we see in the last tast table in the results section that for our effort method XXXXXXXX produced the best results. did we or did we not achieve the results we sought in the abstract