

# Statistical Analysis of Sign Language Recognition with Simple CNN Models

Julia Xu, 293857<sup>1</sup>

<sup>1</sup>Supervisor: Maciej Huk, Ph.D

<sup>2</sup>Institute: Wroclaw University of Science and Technology - Faculty of Information and Communication Technology

<sup>3</sup>Lecture: Advanced Topics in Artificial Intelligence - Project

January 2026

## Contents

<b>1</b>	<b>Research</b>	<b>2</b>
1.1	Research plan . . . . .	2
1.2	Datasets . . . . .	3
1.3	Research environment . . . . .	3
1.3.1	Statistical Analysis . . . . .	4
1.3.2	Power Analysis . . . . .	4
1.4	Measures . . . . .	4
1.4.1	Model Performance Metrics . . . . .	5
1.4.2	Statistical Tests . . . . .	5
1.4.3	Effect sizes . . . . .	5
1.5	Experiment Methodology . . . . .	6
1.5.1	BaselineCNN . . . . .	6
1.5.2	DeeperCNN . . . . .	6
1.5.3	RegularizedCNN . . . . .	7
<b>2</b>	<b>Results</b>	<b>7</b>
2.1	Training Performance . . . . .	7
2.2	Statistical Analysis . . . . .	8
2.2.1	Friedman ANOVA . . . . .	8
2.2.2	Wilcoxon . . . . .	8
2.3	Power Analysis . . . . .	9
<b>3</b>	<b>Summary of research</b>	<b>9</b>

# 1 Research

The purpose of this paper is to perform a statistical and power analysis of a selected problem related to Artificial Intelligence (AI). The chosen task is to perform Sign Language Recognition (SLR) in form of image classification. Further details about the experiment(s) are explained in subsection 1.1. The implementation is published and visible on GitHub.<sup>1</sup>

The motivation behind this project is to analyze whether observed model performance differences occur due real and systematic given a same baseline rather than due to randomness be it from initialization or splits. Naturally, it is to proof whether the training of model variants result in significantly different outcomes when running them as independently and randomized as possible. Hence, the null hypothesis  $H_0$  here is that there is no difference in performance among three selected models. Rejection of this hypothesis would indicate at least one model performs differently. This allows to verify whether the performance differences of three models are statistically significant.

To provide as much randomness as possible, the seed management is conducted by producing MD5-hashed values first which deal as basis for creating randomized seeds used for training and splitting purposes.

Furthermore, Demšar’s paper in 2006 documents thoroughly how to cleanly perform a statistical and power analysis by providing multiple methods to compare various models [1]. This authors recommendations are taken into consideration to construct the research plan for the statistical analysis of SLR through image classification with simple Convolutional Neural Networks (CNNs).

It is assumed that DeeperCNN or RegularizedCNN result the best accuracies while all models flag significant differences. This assumption is researched and explained in the following sections.

## 1.1 Research plan

Figure 1 visualizes the experimental pipeline. There is a chronological nature integrated from left to right. There are three model variants defined: a two-layered CNN (BaselineCNN), a three-layered CNN (DeeperCNN), and a three-layered CNN with a higher regularization factor (RegularizedCNN). All three models use dropout to introduce some kind of regularization while the latter model increases the dropout probability to 0.5 in contrast to 0.3 for the other two models. All three model architectures are visualized in sections 1.5.1, 1.5.2 and section 1.5.3. The activation function is a simple ReLU and max pooling and adapted average pooling are applied in all model variants.

Then the training pipeline is defined with a 5-5-3-5-5-principle: 5 global seeds for training initialization, 5 training epochs, 3 model variants, and 5 times repeated 5-fold CV per seed and model. This results in 25 measurements (5 repetitions for each of the 5 seeds) per variant. To conform to the mentioned 5-5-3-5-5-principle, the training is run with five epochs per fold while doing a 5-fold CV repeated five times per seed.

To evaluate the results, metrics such as accuracy, macro F1-score, and training duration are collected.

These 25 metric values per model are taken to perform the statistical analysis on. First, Friedman ANOVA is done to check whether the experiment results show any significant difference overall. Then, the significant pairs are taken and run through the Wilcoxon test to realize between which pairs does a significant difference exist. Finally, the p-values are run through Hommel for post-hoc correction.

The power analysis follows after the statistical one to measure how big the significant difference between models is and in what directions they go. This is done through running Friedman ANOVA

---

<sup>1</sup><https://github.com/xuju-dev/sign-language-recognition.git>

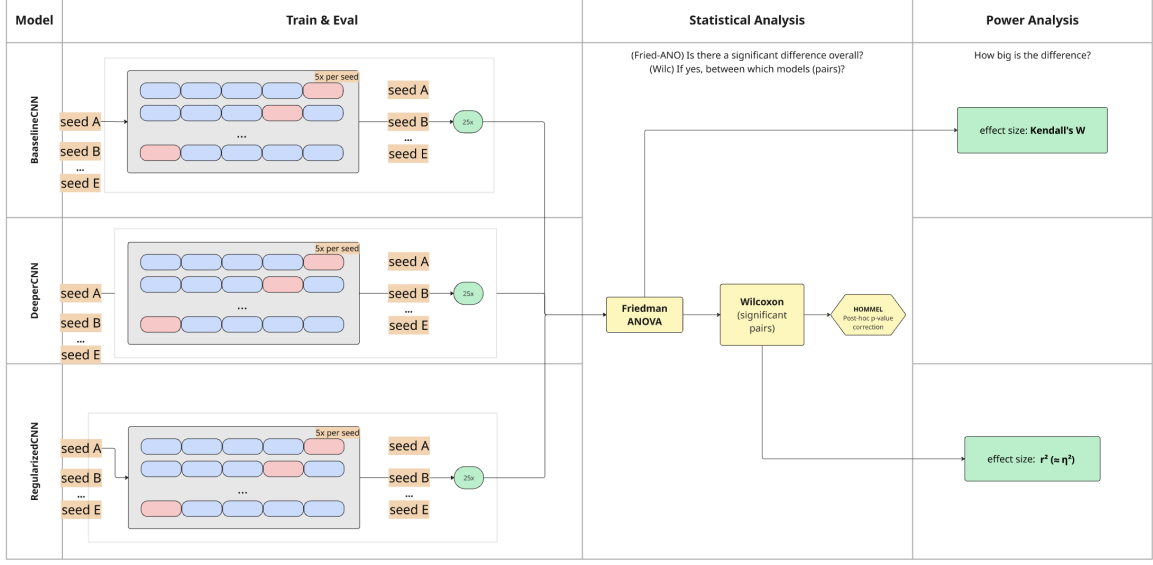


Figure 1: Flow chart of experiment pipeline. Data preprocessing is not included and assumed done already.

results through Kendall's W and mean rank estimator while the Wilcoxon test results are the inputs for calculating  $\eta^2$  for effect size and estimated through Hodges-Lehmann.

## 1.2 Datasets

The chosen dataset for this task of SLR contains in total 87,000 training images in 29 classes [6]. The classes are defined as the alphabetic letter with special characters such as SPACE, DELETE, and NOTHING. Each image contains a hand signature signing one specific character.

A subset of 10% of this dataset is taken to perform the experiments on. From that, the data is merged into one development set (training and validation set) with the original test set with size of 29. The images are originally in 200x200 format and preprocessed to have an image size of 224 to conform to our first draft of experimental pipeline which was with a pretrained model (MobileNet v3 [4, 10, 5]). This route was soon discarded once early training results showed that the pre-trained model was too strong for this simple task. Nonetheless, the same preprocessing pipeline can be used for the second approach with the simpler own constructed CNNs. No data augmentation is added after confirming that early training rounds of the simpler CNNs were showing accuracies low enough. This is further discussed in the later section when we evaluate and discuss the training results.

Although the original dataset contains 29 classes, all experiments were conducted using a 28-class label space due to a configuration oversight. All models were trained and evaluated under identical conditions ensuring the validity of this analysis.

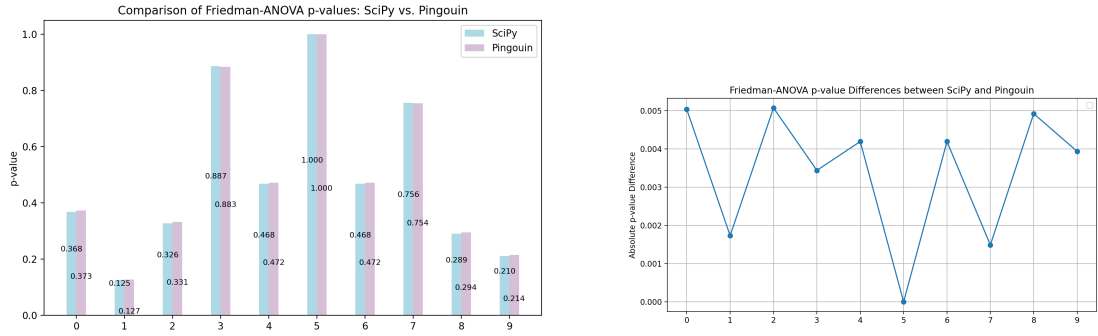
## 1.3 Research environment

The training is run on a Macbook Air with Apple M2 chip, 8-core CPU and 8-core GPU. This Apple hardware does not contain any CUDA supported GPU and provides its own version called

'mps' which is used for this research's experiments. General implementation is done through the programming language Python 3.11.14 with classic libraries such as but not limited to `torch`, `matplotlib`, `numpy`, `scipy`, `pingouin` and `pandas`.

### 1.3.1 Statistical Analysis

The statistical analysis is implemented with a Python library called `scipy`. Before deciding on using `scipy`, a comparison of libraries is conducted to see whether different implementations result in different outcomes given a same simple data input. The comparison was done between the libraries `scipy` and `pingouin`. For Friedman test three vectors are generated with 25 random values. The test is run ten times and the results are shown in Figure 2.



(a) Overview of test results. The x-axis corresponds to the numbered test repetition.

(b) Absolute p-value differences. The x-axis corresponds to the numbered test repetition.

Figure 2: Comparison of statistical test results and absolute p-value differences of ten times repeated statistical test implementation from `scipy` and `pingouin` with three example 1x25 vectors.

Figure 2a shows the overall test results. The p-value is plotted on the y-axis and the repeated test runs are numbered on the x-axis. The `pingouin` implementation of Friedman-ANOVA computes a higher or identical p-value for 80% of the test runs where two out of ten runs resulted in higher p-values for the `scipy` implementation. The numeric differences are visualized in Figure 2b as well where the absolute difference is better highlighted. The highest difference was 0.005 in p-value while the lowest was 0 for test run 5.

The `statsmodels` library is used for conducting post-hoc Hommel correction.

### 1.3.2 Power Analysis

For performing power analysis on our task, manual Python implementation is used. No other established library was found such that the comparison of different implementations is omitted here. The effect sizes are calculated through Kendall's  $W/\eta^2$  and  $r$ . The direction of the differences computed through estimators are omitted in this research due to time constraints.

## 1.4 Measures

This section gives a quick overview and definition of the used metrics in this research. There are metrics to measure model performances and formulas for statistical tests and their effect size.

#### 1.4.1 Model Performance Metrics

Accuracy (Eq. (1)) is a general measure for how good a prediction is based on true positives ( $TP$ ) and true negatives ( $TN$ ) in relation to the full set of samples where  $FP$  denotes false positives, and  $FN$  false negatives [7].

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}, \quad (1)$$

F1 score (macro) is a multi-class based metric using precision and recall to measure F1 score for each class independently and averaging them to equally weight each class independent on their class representation in a given dataset.

$$F1_{\text{macro}} = \frac{1}{C} \sum_{c=1}^C \frac{2 \text{Precision}_c \cdot \text{Recall}_c}{\text{Precision}_c + \text{Recall}_c}. \quad (2)$$

The training duration is logged in minutes and records the time per fold. This time is then added up to conclude the duration for a full run of 5-fold CV (here may be also referred as 'repeat/repetition') which is taken to compare with other experiment groups.

#### 1.4.2 Statistical Tests

The statistical analysis consists of Friedman ANOVA test and Wilcoxon signed-ranks test [8, 9]. Friedman ANOVA is used to test on more than two models. Once the Friedman ANOVA test returns a significant difference in the compared models the Wilcoxon signed-ranks test is run to compare two specific models to explore which model/group pair is statistically significant. After that, Hommel is applied to conduct a post-hoc p-value correction [3].

#### 1.4.3 Effect sizes

To determine the power of the statistical tests, the effect sizes are calculated. For Friedman ANOVA Kendall's W (Eq. (3)) is used where  $\chi^2$  is the Friedman ANOVA test statistics,  $n$  the number of conducted values per metric, and  $k$  the number of models/conditions.

The effect size  $W_K$ <sup>2</sup> is considered large when  $\geq 0.5$ , medium when  $\geq 0.3$ , and small when  $\geq 0.1$  [2].

$$W_K = \frac{\chi^2}{n(k-1)} \quad (3)$$

For Wilcoxon test the equations (4)-(6) are used. Eq. (4) extracts the number of non-zero differences, Eq. (5) computes the expected value and standard deviation of the observed Wilcoxon signed-rank statistics  $W$ . Finally, the Z-score is approximated and used to compute the effect size  $r$  (Eq. (4)) where  $N$  is the total number of paired observations.

$$n = \sum_{i=1}^N \mathbf{1}(x_i - y_i \neq 0) \quad (4)$$

$$\mu_W = \frac{n(n+1)}{4}, \sigma_W = \sqrt{\frac{n(n+1)(2n+1)}{24}} \quad (5)$$

---

<sup>2</sup>Here  $W_K$  is referred to Kendall's W to avoid confusion with Wilcoxon's test statistics  $W$

$$Z = \frac{W - \mu_W}{\sigma_W}, r = \frac{|Z|}{\sqrt{N}} \quad (6)$$

The effect size  $r$  has the same interpretation thresholds as for  $W_K$  [2].

## 1.5 Experiment Methodology

The three experimental groups are equivalent to the model variants: BaselineCNN, DeeperCNN, RegularizedCNN. All experimental groups conduct the same training pipeline with the same pre-processed dataset. Each group gets the same five seeds. The difference in these experiments is limited to the used model variant to ensure a stable environment for the upcoming statistical analysis.

The statistical analysis contains a Friedman ANOVA test which is then followed by pairwise Wilcoxon test. The final layer provides a post-hoc p-value correction through Hommel.

Finally, the effect sizes are computed for the given statistical tests to show which model comparisons show the highest power in significance.

### 1.5.1 BaselineCNN

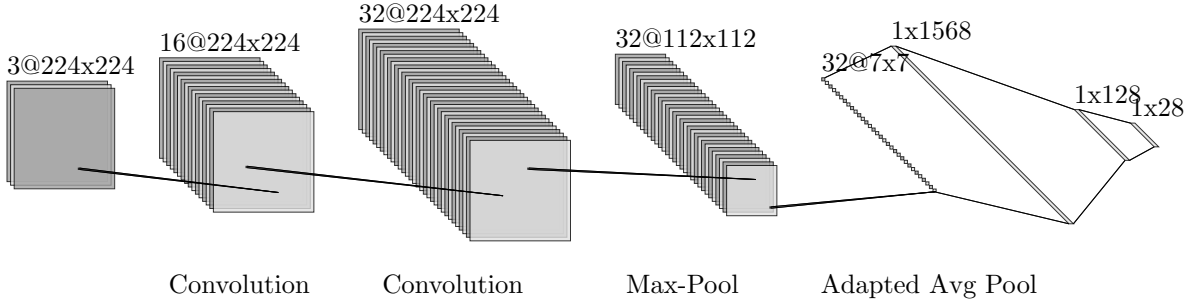


Figure 3: Model architecture of BaselineCNN illustrating size changes throughout the pipeline. Activation and dropout layers are omitted since they do not change any sizes. At the most left is the input image of size 224x224. Convolutional layers have a kernel size of 3 and padding of 1, hence not changing the spatial size. Max-Pool layers have size 2 and stride 2 and the average adapted pooling layer compresses its input into a 7x7 spatial size.

BaselineCNN consists of a 2-layered CNN with one layer of max pooling and a second resizing layer using adapted average pooling. All activations are using ReLU. All convolutional kernels have size 3 with padding of 1 and the max pooling layers half the spatial size with the last (adapted) pooling layer reducing the spatial space to 7x7. A dropout layer is applied once at the end before flattening and has the probability of 0.3.

Figure 3 shows the model architecture of BaselineCNN and illustrates how each layer reduces the spatial size.

### 1.5.2 DeeperCNN

DeeperCNN takes the beginning of BaselineCNN and includes an additional convolutional layer with a following max pooling layer before the adapted average pooling layer ending with an adapted average pooling space of 64@7x7. The dropout probability stays the same at 0.3 and the filter sizes and padding values stay consistent with the simpler CNN model.



Figure 4: Full model architecture pipeline of BaselineCNN. This model has 2 convolutional layers and 1 maximum pooling layer. Convolutions have kernel size 3 and padding 1. Max pooling layers have a filter of size 2 with stride 2 and the adapted average pooling reduces the spatial space to 7x7. Dropout probability is 0.3.



Figure 5: Full model architecture pipeline of DeeperCNN. This model has 3 convolutional layers and 2 maximum pooling layers. Convolutions have kernel size 3 and padding 1. Max pooling layers have a filter of size 2 with stride 2 and the adapted average pooling reduces the spatial space to 7x7. Dropout probability is 0.3.

### 1.5.3 RegularizedCNN

RegularizedCNN has the same architecture as DeeperCNN but with an increased dropout probability of 0.5 elevating the regularization factor in this model.



Figure 6: Full model architecture pipeline of RegularizedCNN. Convolutions have kernel size 3 and padding 1. Max pooling layers have a filter of size 2 with stride 2 and the adapted average pooling reduces the spatial space to 7x7. Dropout probability is 0.5.

## 2 Results

### 2.1 Training Performance

As shown in Figure 7, DeeperCNN results higher accuracy values and F1-scores compared to the other two models. However, DeeperCNN also takes the longest to train out of these three models. BaselineCNN shows the shortest training duration which conforms to expected outcomes since it is the smallest model out of these three. RegularizedCNN shows shorter but similar training durations compared to DeeperCNN but the accuracy values and F1-scores underperform the base model. This shows that added regularization such as higher dropout probability does not increase a model's performance in the task of SLR with out given dataset. On the contrary, a simpler model such as BaselineCNN results in better metrics regarding correct classifications.

The model with an added convolutional layer and 0.3 dropout rate shows the best overall performance in terms of model performance metrics.

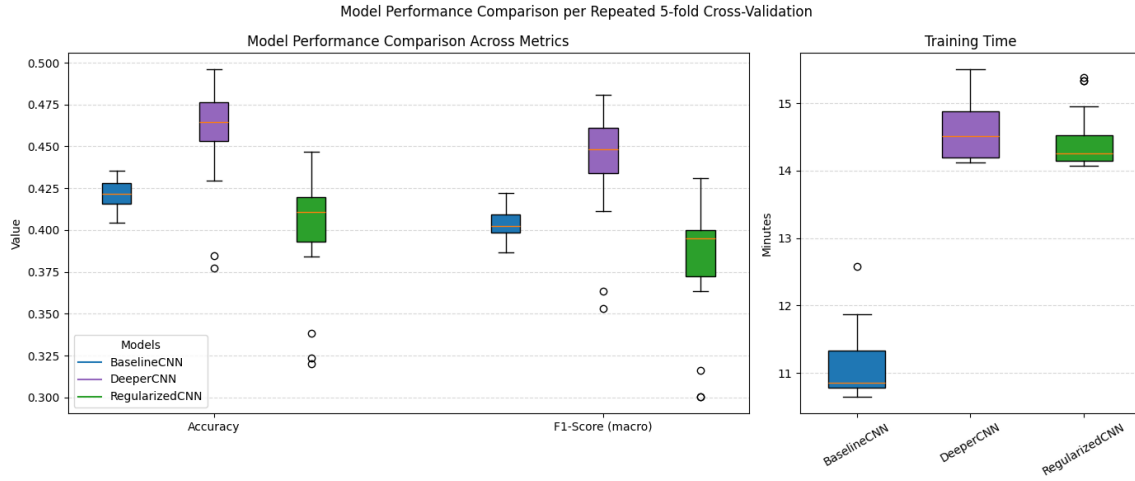


Figure 7: Overall model comparison per repeated 5-fold CV. In total there are 25 values per model (5 seeds times 5 repetitions) [Values add up to 1].

Next, the statistical analysis indicates whether these models are statistically significant and if so which model pairs are significantly different and how likely so.

## 2.2 Statistical Analysis

### 2.2.1 Friedman ANOVA

The Friedman-ANOVA test shows that the models are statistically significant in both metrics accuracy as well as F1-score. Both p-values are well below the set  $\alpha$ -level of 5% where the metric accuracy has a lower value.

Table 1: Friedman ANOVA test p-values.

Metric	p-value	Significant ( $\alpha = 0.05$ )
Accuracy	$6.96 \times 10^{-8}$	✓
F1-score (macro)	$1.17 \times 10^{-7}$	✓

### 2.2.2 Wilcoxon

The Wilcoxon paired signed-rank test result show that a significant difference can be found on all model pairs given the same set  $\alpha$ -level of 0.05. The highest p-value results from comparing BaselineCNN with RegularizedCNN. The post-hoc Hommel correction shows the biggest effect in the BaselineCNN-DeeperCNN model pair where the correction increased the Wilcoxon p-value by  $3.20 \times 10^{-3}$ . Interestingly, BaselineCNN-RegularizedCNN resulted the same p-value even after Hommel. Nonetheless, these test results show that there exists a statistical significance between all three model pairs and confirms our initial assumption.



Table 2: Wilcoxon signed-rank test with Hommel correction for pairwise comparisons of CNN models (accuracy).

Comparison	Wilcoxon $p$ -value	Hommel-corrected $p$	Significant ( $\alpha = 0.05$ )
BaselineCNN vs DeeperCNN	$3.19 \times 10^{-5}$	$6.39 \times 10^{-5}$	✓
BaselineCNN vs RegularizedCNN	$3.78 \times 10^{-3}$	$3.78 \times 10^{-3}$	✓
DeeperCNN vs RegularizedCNN	$1.19 \times 10^{-7}$	$3.58 \times 10^{-7}$	✓

## 2.3 Power Analysis

The following effect sizes show how powerful the resulted significances are. Table 3 shows the effect size for Friedman ANOVA. Both accuracy and F1-score have a large effect size showing a strong and significant difference in models.

Table 3: Friedman ANOVA effect size (Kendall’s  $W / \eta^2$ ) for evaluation metrics.

Metric	Kendall’s $W (\eta^2)$
Accuracy	0.6592
F1-score	0.6384

Table 4 shows the effect size for the Wilcoxon test. The largest effect size occurs when comparing DeeperCNN with RegularizedCNN while the smallest effect size conforms to the BaselineCNN-RegularizedCNN model pair. Overall, all computed effect sizes here indicate a strong difference in model pairs where the biggest difference is found in higher regularization factor with the same amount of convolutional layers. Interestingly enough, the smallest difference is found between our simplest and most tweaked model with an effect size of 0.5624.

Table 4: Wilcoxon signed-rank effect size ( $r$ ) for pairwise CNN model comparisons (accuracy).

Comparison	Effect size $r$
BaselineCNN vs DeeperCNN	0.7561
BaselineCNN vs RegularizedCNN	0.5624
DeeperCNN vs RegularizedCNN	0.8691

## 3 Summary of research

Overall, this research shows comparing multiple models need a certain statistical pipeline. Here, three models were compared to show if a significant difference exists at all. The Friedman ANOVA test shows that there indeed exists a statistical significance.

Furthermore, the Wilcoxon signed rank test is applied to evaluate which model pairs conduct a significant difference and all model pair permutations achieved a  $p$ -value that confirms a statistical significance between them.

Finally, the effect size of both statistical tests show that the extracted significance is to be seen with large importance. While accuracy and F1-score have similar effect sizes when comparing all

three models the power of the paired tests shows that BaselineCNN and DeeperCNN have the largest importance regarding statistic significance. This confirms our initial assumptions that all three models are statistically different.

To sum up, DeeperCNN performs the best in terms of simple model performance metrics and shows significant difference to all other model variants with the biggest difference with Regularized-CNN.

## References

- [1] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *J. Mach. Learn. Res.*, 7:1–30, December 2006.
- [2] Fernanda Fiel Peres. Effect sizes for nonparametric tests. *Biochemia Medica*, 36(1), December 2025.
- [3] G. HOMMEL. A stagewise rejective multiple test procedure based on a modified bonferroni test. *Biometrika*, 75(2):383–386, 1988.
- [4] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1314–1324, 2019.
- [5] Hugging Face. Model card for mobilenetv3\_small\_100.lamb\_in1k. [https://huggingface.co/timm/mobilenetv3\\_small\\_100.lamb\\_in1k](https://huggingface.co/timm/mobilenetv3_small_100.lamb_in1k). [Accessed 21-11-2025].
- [6] Khan, Anas and Nagaraj, Akash. Asl alphabet. <https://www.kaggle.com/datasets/grassknoted/asl-alphabet>. [Accessed 03-11-2025].
- [7] Oona Rainio, Jarmo Teuho, and Riku Klén. Evaluation metrics and statistical tests for machine learning. *Scientific Reports*, 14(1), March 2024.
- [8] scipy. friedmanchisquare SciPy v1.16.2 Manual. <https://docs.scipy.org/doc/scipy-1.16.2/reference/generated/scipy.stats.friedmanchisquare.html>. [Accessed 27-11-2025].
- [9] scipy. wilcoxonSciPy v1.16.2 Manual. <https://docs.scipy.org/doc/scipy-1.16.2/reference/generated/scipy.stats.wilcoxon.html>. [Accessed 27-11-2025].
- [10] Ross Wightman. Pytorch image models. <https://github.com/huggingface/pytorch-image-models>, 2019.