# LSTM Fully Convolutional Networks for Time Series Classification

Fazle Karim[1*], Somshubra Majumdar[2*], Houshang Darabi[1], *Senior Member, IEEE, and Shun Chen*[1]

*Abstract*—**Fully convolutional neural networks (FCN) have been shown to achieve state-of-the-art performance on the task of classifying time series sequences. We propose the augmentation of fully convolutional networks with long short term memory recurrent neural network (LSTM RNN) sub-modules for time series classification. Our proposed models significantly enhance the performance of fully convolutional networks with a nominal increase in model size and require minimal preprocessing of the dataset. The proposed Long Short Term Memory Fully Convolutional Network (LSTM-FCN) achieves state-of-the-art performance compared to others. We also explore the usage of attention mechanism to improve time series classification with the Attention Long Short Term Memory Fully Convolutional Network (ALSTM-FCN). Utilization of the attention mechanism allows one to visualize the decision process of the LSTM cell. Furthermore, we propose fine-tuning as a method to enhance the performance of trained models. An overall analysis of the performance of our model is provided and compared to other techniques.**

*Keywords*—*Convolutional Neural Network, Long Short Term Memory Recurrent Neural Network, Time Series Classification*

## I. INTRODUCTION

Over the past decade, there has been an increased interest in time series classification. Time series data is ubiquitous, existing in weather readings, financial recordings, industrial observations, and psychological signals [1]. In this paper two deep learning models to classify time series datasets are proposed, both of which outperform existing state-of-the-art models.

A plethora of research have been done using feature-based approaches or methods to extract a set of features that represent time series patterns. Bag-of-Words (BoW) [2], Bag-of-features (TSBF) [3], Bag-of-SFA-Symbols (BOSS) [4], BOSSVS [5], Word ExtrAction for time Series cLassification (WEASEL) [6], have obtained promising results in the field. Bag-of-words quantizes the extracted features and feeds the BoW into a classifier. TSBF extracts multiple subsequences of random local information, which a supervised learner condenses into a cookbook used to predict time series labels. BOSS introduces a combination of a distance based classifier and histograms. The histograms represent substructures of a time series that are created using a symbolic Fourier approximation. BOSSVS extends this method by proposing a vector space model to reduce time complexity while maintaining performance.

WEASEL converts time series into feature vectors using a sliding window. Machine learning algorithms utilize these feature vectors to detect and classify the time series. All these classifiers require heavy feature extraction and feature engineering.

Ensemble algorithms also yield state-of-the-art performance with time series classification problems. Three of the most successful ensemble algorithms that integrate various features of a time series are Elastic Ensemble (PROP) [7], a model that integrates 11 time series classifiers using a weighted ensemble method, Shapelet ensemble (SE) [8], a model that applies a heterogeneous ensemble onto transformed shapelets, and a flat collective of transform based ensembles (COTE) [8], a model that fuses 35 various classifiers into a single classifier.

Recently, deep neural networks have been employed for time series classification tasks. Multi-scale convolutional neural network (MCNN) [9], fully convolutional network (FCN) [10], and residual network (ResNet) [10] are deep learning approaches that take advantage of convolutional neural networks (CNN) for end-to-end classification of univariate time series. MCNN uses down-sampling, skip sampling and sliding window to preprocess the data. The performance of the MCNN classifier is highly dependent on the preprocessing applied to the dataset and the tuning of a large set of hyperparameters of that model. On the other hand, FCN and ResNet do not require any heavy preprocessing on the data or feature engineering. In this paper, we improve the performance of FCN by augmenting the FCN module with either a Long Short Term Recurrent Neural Network (LSTM RNN) sub-module , called LSTM-FCN, or a LSTM RNN with attention, called ALSTM-FCN. Similar to FCN, both proposed models can be used to visualize the Class Activation Maps (CAM) of the convolutional layers to detect regions that contribute to the class label. In addition, the Attention LSTM can also be used detect regions of the input sequence that contribute to the class label through the context vector of the Attention LSTM cells. A major advantages of the LSTM-FCN and ALSTM-FCN models is that it does not require heavy preprocessing or feature engineering. Results indicate the new proposed models, LSTM-FCN and ALSTM-FCN, dramatically improve performance on the University of California Riverside (UCR) Benchmark datasets [11]. LSTM-FCN and ALSTM-FCN produce better results than several state-of-the-art ensemble algorithms on a majority of the UCR Benchmark datasets.

This paper proposes two deep learning models for end-to-end time series classification. The proposed models do not require heavy preprocessing on the data or feature engineering. Both the models are tested on all 85 UCR time series benchmarks and outperform most of the state-of-the-art

*Equal contribution.

[1]Mechanical and Industrial Engineering, University of Illinois at Chicago, Chicago,IL

[2]Computer Science, University of Illinois at Chicago, Chicago, IL

models. The remainder of the paper is organized as follows. Section II reviews the background work. Section III presents the architecture of the proposed models. Section IV analyzes and discusses the experiments performed. Finally, conclusions are drawn in Section V.

## II. BACKGROUND WORKS

### A. Temporal Convolutions

The input to a Temporal Convolutional Network is generally a time series signal. As stated in *Lea et al.* [12], let $X_t \in \mathbb{R}^{F_0}$ be the input feature vector of length $F_0$ for time step $t$ for $0 < t \le T$. Note that the time T may vary for each sequence, and we denote the number of time steps in each layer as $T_l$. The true action label for each frame is given by $y_t \in \{1, ..., C\}$, where C is the number of classes.

Consider $L$ convolutional layers. We apply a set of 1D filters on each of these layers that capture how the input signals evolve over the course of an action. According to *Lea et al.* [12], the filters for each layer are parameterized by tensor $W^{(l)} \in \mathbb{R}^{F_l \times d \times F_{l-1}}$ and biases $b^{(l)} \in \mathbb{R}^{F_l}$, where $l \in \{1, ..., L\}$ is the layer index and $d$ is the filter duration. For the $l$-th layer, the $i$-th component of the (unnormalized) activation $\hat{\mathbf{E}}_t^{(l)} \in \mathbb{R}^{F_l}$ is a function of the incoming (normalized) activation matrix $E^{(l-1)} \in \mathbb{R}^{F_{l-1} \times T_{l-1}}$ from the previous layer

$$\hat{\mathbf{E}}_{i,t}^{(l)} = f\left(b_i^{(l)} + \sum_{t'=1}^{d}\left\langle W_{i,t',.}^{(l)}, E_{.,t+d-t'}^{(l-1)}\right\rangle\right) \quad (1)$$

for each time $t$ where $f(\cdot)$ is a Rectified Linear Unit.

We use Temporal Convolutional Networks as a feature extraction module in a Fully Convolutional Network (FCN) branch. A basic convolution block consists of a convolution layer, followed by batch normalization [13], followed by an activation function, which can be either a Rectified Linear Unit or a Parametric Rectified Linear Unit [14].

### B. Recurrent Neural Networks

Recurrent Neural Networks, often shortened to RNNs, are a class of neural networks which exhibit temporal behaviour due to directed connections between units of an individual layer. As reported by *Pascanu et al.* [15], recurrent neural networks maintain a hidden vector $\mathbf{h}$, which is updated at time step $t$ as follows:

$$\mathbf{h}_t = \tanh(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{I}\mathbf{x}_t), \quad (2)$$

tanh is the hyperbolic tangent function, $\mathbf{W}$ is the recurrent weight matrix and $\mathbf{I}$ is a projection matrix. The hidden state $\mathbf{h}$ is used to make a prediction

$$\mathbf{y}_t = \text{softmax}(\mathbf{W}\mathbf{h}_{t-1}), \quad (3)$$

softmax provides a normalized probability distribution over the possible classes, $\sigma$ is the logistic sigmoid function and $\mathbf{W}$ is a weight matrix. By using $\mathbf{h}$ as the input to another RNN, we can stack RNNs, creating deeper architectures

$$\mathbf{h}_t^l = \sigma(\mathbf{W}\mathbf{h}_{t-1}^l + \mathbf{I}\mathbf{h}_t^{l-1}). \quad (4)$$

### C. Long Short-Term Memory RNNs

Long short-term memory recurrent neural networks are an improvement over the general recurrent neural networks, which possess a vanishing gradient problem. As stated in *Hochreiter et al.* [16], LSTM RNNs address the vanishing gradient problem commonly found in ordinary recurrent neural networks by incorporating gating functions into their state dynamics. At each time step, an LSTM maintains a hidden vector $\mathbf{h}$ and a memory vector $\mathbf{m}$ responsible for controlling state updates and outputs. More concretely, *Graves et al.* [17] define the computation at time step $t$ as follows :

$$
\begin{aligned}
\mathbf{g}^u &= \sigma(\mathbf{W}^u\mathbf{h}_{t-1} + \mathbf{I}^u\mathbf{x}_t) \\
\mathbf{g}^f &= \sigma(\mathbf{W}^f\mathbf{h}_{t-1} + \mathbf{I}^f\mathbf{x}_t) \\
\mathbf{g}^o &= \sigma(\mathbf{W}^o\mathbf{h}_{t-1} + \mathbf{I}^o\mathbf{x}_t) \\
\mathbf{g}^c &= \tanh(\mathbf{W}^c\mathbf{h}_{t-1} + \mathbf{I}^c\mathbf{x}_t) \\
\mathbf{m}_t &= \mathbf{g}^f \odot \mathbf{m}_{t-1} + \mathbf{g}^u \odot \mathbf{g}^c \\
\mathbf{h}_t &= \tanh(\mathbf{g}^o \odot \mathbf{m}_t)
\end{aligned}
\quad (5)
$$

where $\sigma$ is the logistic sigmoid function, $\odot$ represents elementwise multiplication, $\mathbf{W}^u, \mathbf{W}^f, \mathbf{W}^o, \mathbf{W}^c$ are recurrent weight matrices and $\mathbf{I}^u, \mathbf{I}^f, \mathbf{I}^o, \mathbf{I}^c$ are projection matrices.

While LSTMs possess the ability to learn temporal dependencies in sequences, they have difficulty with long term dependencies in long sequences. The attention mechanism proposed by *Bahdanau et al.* [18] can help the LSTM RNN learn these dependencies.

### D. Attention Mechanism

The attention mechanism is a technique often used in neural translation of text, where a context vector $C$ is conditioned on the target sequence $y$. As discussed in *Bahdanau et al.* [18], the context vector $c_i$ depends on a sequence of *annotations* $(h_1, ..., h_{T_x})$ to which an encoder maps the input sequence. Each annotation $h_i$ contains information about the whole input sequence with a strong focus on the parts surrounding the $i$-th word of the input sequence.

The context vector $c_i$ is then computed as a weighted sum of these annotations $h_i$:

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j. \quad (6)$$

The weight $\alpha_{ij}$ of each annotation $h_j$ is computed by :

$$\alpha_{ij} = \frac{exp(e_{ij})}{\sum_{k=1}^{T_x} exp(e_{ik})} \quad (7)$$

where $e_{ij} = a(s_{i-1}, h_j)$ is an *alignment model*, which scores how well the input around position $j$ and the output at position $i$ match. The score is based on the RNN hidden state $s_{i1}$ and the $j$-th annotation $h_j$ of the input sentence.

*Bahdanau et al.* [18] parametrize the alignment model $a$ as a feedforward neural network which is jointly trained with all the other components of the model. The alignment model directly computes a soft alignment, which allows the gradient of the cost function to be backpropagated.
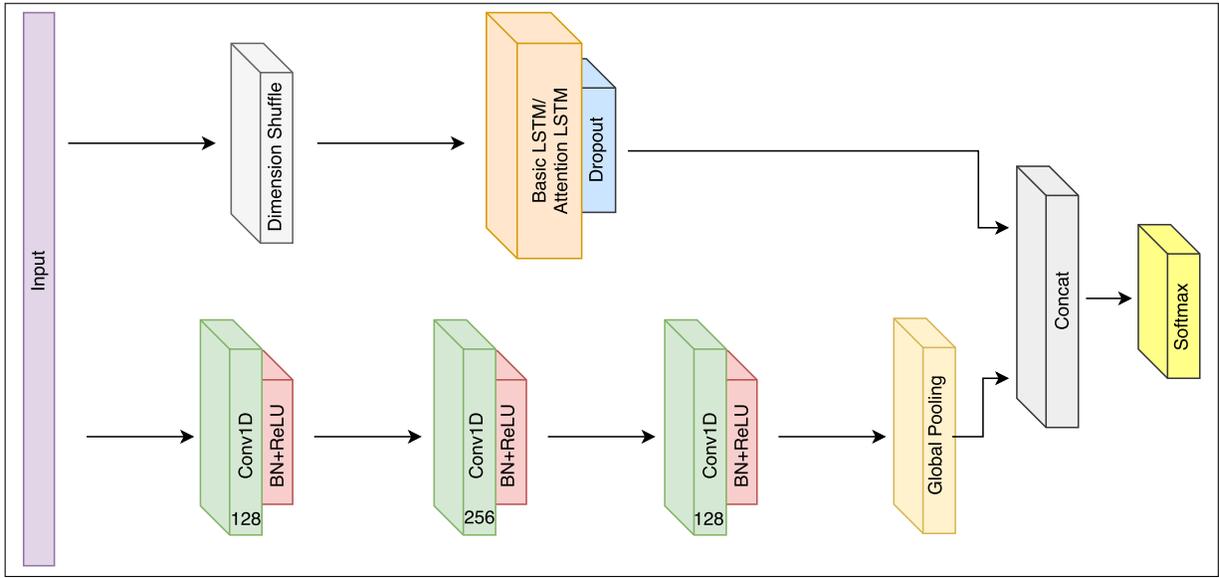
Fig. 1: The LSTM-FCN architecture. LSTM cells can be replaced by Attention LSTM cells to construct the ALSTM-FCN architecture.

## III. LSTM FULLY CONVOLUTIONAL NETWORK

### A. Network Architecture

Temporal convolutions have proven to be an effective learning model for time series classification problems [10]. Fully Convolutional Networks comprised of temporal convolutions are typically used as feature extractors, and global average pooling [19] is used to reduce the number of parameters in the model prior to classification. In the proposed models, the fully convolutional block is augmented by an LSTM block followed by dropout [20], as shown in Fig.1.

The fully convolutional block consists of three stacked temporal convolutional blocks with filter sizes of 128, 256, and 128 respectively. Each convolutional block is identical to the convolution block in the CNN architecture proposed by *Wang et al.* [10]. Each block consists of a temporal convolutional layer, which is accompanied by batch normalization [13] (momentum of 0.99, epsilon of 0.001) followed by a ReLU activation function. Finally, global average pooling is applied following the final convolution block.

Simultaneously, the time series input is conveyed into a dimension shuffle layer (explained more in Section III-B). The transformed time series from the dimension shuffle is then passed into the LSTM block. The LSTM block comprises of either a general LSTM layer or an Attention LSTM layer, followed by a dropout. The output of the global pooling layer and the LSTM block is concatenated and passed onto a softmax classification layer.

### B. Network Input

The fully convolutional block and LSTM block perceive the same time series input in two different views. The fully convolutional block views the time series as a univariate time series with multiple time steps. If there is a time series of length $N$, the fully convolutional block will receive the data in $N$ time steps.

Contrarily, the LSTM block in the proposed architecture receives the input time series as a multivariate time series with a single time step. This is accomplished by the dimension shuffle layer, which transposes the temporal dimension of the time series. A univariate time series of length $N$, after transformation, will be viewed as a multivariate time series (having $N$ variables) with a single time step.

This approach is key to the enhanced performance of the proposed architecture. In contrast, when the LSTM block received the univariate time series with $N$ time steps, the performance was significantly reduced due to rapid overfitting on small short-sequence UCR datasets and a failure to learn long term dependencies in the larger long-sequence UCR datasets.

### C. Fine-Tuning of Models

Transfer learning is a technique wherein the knowledge gained from training a model on a dataset can be reused when training the model on another dataset, such that the domain of the new dataset has some similarity with the prior domain [21]. Similarly, *fine-tuning* can be described as transfer learning on the same dataset.

The training procedure can thus be split into two distinct phases. In the *initial phase*, the optimal hyperparameters for the model are selected for a given dataset. The model is then trained on the given dataset with these hyperparameter settings. In the second step, we apply *fine-tuning* to this initial model.
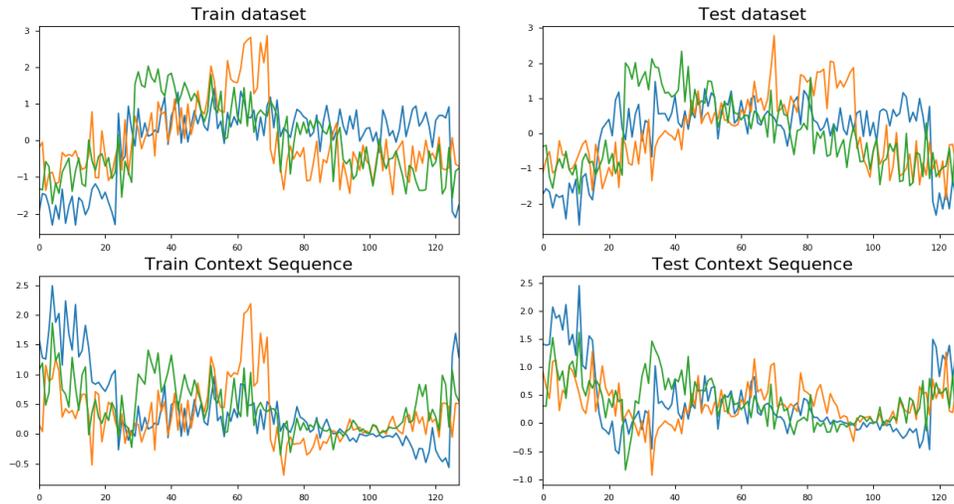
Fig. 2: Visualization of context vector on CBF dataset.

The procedure of transfer learning is iterated over in the *fine-tuning phase*, using the original dataset. Each repetition is initialized using the model weight of the previous iteration. At each iteration the learning rate is halved. Furthermore, the batch size is halved once every alternate iteration. This is done until the initial learning rate is $1e-4$ and batch size is 32. The procedure is repeated $K$ times, where K is an arbitrary constant, generally set as 5.

---

**Algorithm 1** Fine-tuning
---
1: **for** $i < K$ **do**
2:     $model_{weights} \leftarrow initial\_model_{weights}$
3:     *Train(model, initial_lr, batchsize)*
4:     $initial\_model_{weights} \leftarrow model_{weights}$
5:     $i \leftarrow i + 1$
6:     $initial\_lr \leftarrow updateLearningRate(initial\_lr, i)$
7:     $batchsize \leftarrow updateBatchsize(batchsize, i)$

---

## IV. EXPERIMENTS

The proposed models have been tested on all 85 UCR time series datasets [11]. The FCN block was kept constant throughout all experiments. The optimal number of LSTM cells was found by hyperparameter search over a range of 8 cells to 128 cells. The number of training epochs was generally kept constant at 2000 epochs, but was increased for datasets where the algorithm required a longer time to converge. Initial batch size of 128 was used, and halved for each successive iteration of the *fine-tuning* algorithm. A high dropout rate of 80% was used after the LSTM or Attention LSTM layer to combat overfitting. Class imbalance was handled via a class weighing scheme inspired by *King et al.* [22].

All models were trained via the Adam optimizer [23], with an initial learning rate of $1e-3$ and a final learning rate of $1e-4$. All convolution kernels were initialized with the initialization proposed by *He et al.* [24]. The learning rate was reduced by a factor of $1/\sqrt[3]{2}$ every 100 epochs of no improvement in the validation score, until the final learning rate was reached. No additional preprocessing was done on the UCR datasets as they have close to zero mean and unit variance. All models were *fine-tuned*, and scores stated in Table I refer to the scores obtained by models prior to and after *fine-tuning*. [1]

### A. Evaluation Metrics

In this paper, the proposed model was evaluated using accuracy, rank based statistics, and the mean per class error as stated by *Wang et al.* [10].

The rank-based evaluations used are the arithmetic rank, geometric rank, and the Wilcoxon signed rank test. The arithmetic rank is the arithmetic mean of the rank of dataset. The geometric rank is the geometric mean of the rank of each dataset. The Wilcoxson signed rank test is used to compare the median rank of the proposed model and the existing state-of-the-art models. The null hypothesis and alternative hypothesis are as follows:

$$H_o : Median_{proposed\ model} = Median_{state\text{-}of\text{-}the\text{-}art\ model}$$

$$H_a : Median_{proposed\ model} \neq Median_{state\text{-}of\text{-}the\text{-}art\ model}$$

Mean Per Class Error (MPCE) is defined as the arithmetic mean of the per class error (PCE),

$$PCE_k = \frac{1 - accuracy}{number\ of\ unique\ classes}$$

$$MPCE = \frac{1}{K} \sum PCE_K.$$

---

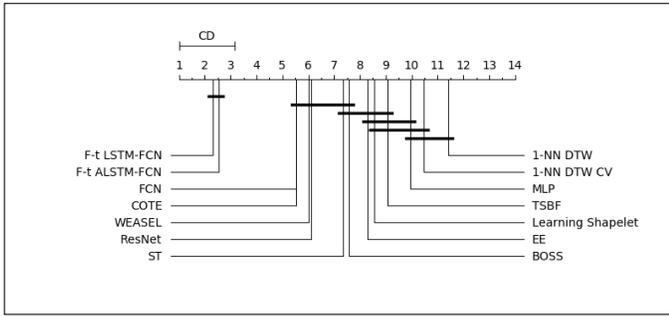[1]The codes and weights of each models are available at https://github.com/houshd/LSTM-FCN

Fig. 3: Critical difference diagram of the arithmetic means of the ranks

## B. Results

Fig. 2 is an example of the visual representation of the Attention LSTM cell on the "CBF" dataset. The points in the figure where the sequences are "squeezed" together are points at which all the classes have the same weight. These are the points in the time series at which the Attention LSTM can correctly identify the class. This is further supported by visual inspection of the actual time series. The squeeze points are points where each of the classes can be distinguished from each other, as shown in Fig. 2.

The performance of the proposed models on the UCR datasets are summarized in Table I. The colored cells are cells that outperform the state-of-the-art model for that dataset. Both proposed models, the ALSTM-FCN model and the LSTM-FCN model, with both phases, without *fine-tuning* (Phase 1) and with *fine-tuning* (Phase 2), outperforms the state-of-the-art models in at least 43 datasets. The average arithmetic rank in Fig. 3 indicates the superiority of our proposed models over the existing state-of-the-art models. This is further validated using the Wilcoxon signed rank test, where the p-value of each of the proposed models are less than 0.05 when compared to existing state-of-the-art models, Table II.

The Wilcoxon Signed Test also provides evidence that *fine-tuning* maintains or improves the overall accuracy on each of the proposed models. The MPCE of the LSTM-FCN and ALSTM-FCN models was found to reduce by 0.0035 and 0.0007 respectively when *fine-tuning* was applied. *Fine-tuning* improves the accuracy of the LSTM-FCN models on a greater number of datasets as compared to the ALSTM-FCN models. We postulate that this discrepancy is due to the fact that the LSTM-FCN model contains fewer total parameters than the ALSTM-FCN model. This indicates a lower rate of overfitting on the UCR datasets. As a consequence, *fine-tuning* is more effective on the LSTM-FCN models for the UCR datasets.

A significant drawback of *fine-tuning* is that it requires more training time due to the added computational complexity of re-training the model using smaller batch sizes. The disadvantages of *fine-tuning* are mitigated when using the ALSTM-FCN within Phase 1. At the end of Phase 1, the ALSTM-FCN

[2]Green cells designate instances where our performance matches or exceeds state-of-the-art results. ∗ denotes model with best performance.

TABLE I: Performance comparison of proposed models with the rest.

| Dataset | Existing SOTA [6, 10] | LSTM-FCN | F-t LSTM-FCN | ALSTM-FCN | F-t ALSTM-FCN |
|---|---|---|---|---|---|
| Adiac | 0.8570 | 0.8593 | 0.8849 | 0.8670 | 0.8900* |
| ArrowHead | 0.8800 | 0.9086 | 0.9029 | 0.9257* | 0.9200 |
| Beef | 0.9000 | 0.9000 | 0.9330 | 0.9333* | 0.9333* |
| BeetleFly | 0.9500 | 0.9500 | 1.0000* | 1.0000* | 1.0000* |
| BirdChicken | 0.9500 | 1.0000* | 1.0000* | 1.0000* | 1.0000* |
| Car | 0.9330 | 0.9500 | 0.9670 | 0.9667 | 0.9833* |
| CBF | 1.0000 | 0.9978 | 1.0000* | 0.9967 | 0.9967 |
| ChloConc | 0.8720 | 0.8099 | 1.0000* | 0.8070 | 0.8070 |
| CinC_ECG | 0.9949 | 0.8862 | 0.9094 | 0.9058 | 0.9058 |
| Coffee | 1.0000 | 1.0000* | 1.0000* | 1.0000* | 1.0000* |
| Computers | 0.8480 | 0.8600 | 0.8600 | 0.8640* | 0.8640* |
| Cricket_X | 0.8210 | 0.8077 | 0.8256* | 0.8051 | 0.8051 |
| Cricket_Y | 0.8256 | 0.8179 | 0.8256* | 0.8205 | 0.8205 |
| Cricket_Z | 0.8154 | 0.8103 | 0.8257 | 0.8308 | 0.8333* |
| DiaSizeRed | 0.9670 | 0.9673 | 0.9771* | 0.9739 | 0.9739 |
| DistPhxAgeGp | 0.8350 | 0.8600 | 0.8600 | 0.8625* | 0.8600 |
| DistPhxCorr | 0.8200 | 0.8250 | 0.8217 | 0.8417* | 0.8383 |
| DistPhxTW | 0.7900 | 0.8175 | 0.8100 | 0.8175 | 0.8200* |
| Earthquakes | 0.8010 | 0.8354* | 0.8261 | 0.8292 | 0.8292 |
| ECG200 | 0.9200 | 0.9000 | 0.9200* | 0.9100 | 0.9200 |
| ECG5000 | 0.9482 | 0.9473 | 0.9478 | 0.9484* | 0.9496* |
| ECGFiveDays | 1.0000 | 0.9919 | 0.9942 | 0.9954 | 0.9954 |
| ElectricDevices | 0.7993 | 0.7681 | 0.7633 | 0.7672 | 0.7672 |
| FaceAll | 0.9290 | 0.9402 | 0.9680 | 0.9657 | 0.9728* |
| FaceFour | 1.0000 | 0.9432 | 0.9772 | 0.9432 | 0.9432 |
| FacesUCR | 0.9580 | 0.9293 | 0.9898* | 0.9434 | 0.9434 |
| FiftyWords | 0.8198 | 0.8044 | 0.8066 | 0.8242 | 0.8286* |
| Fish | 0.9890 | 0.9829 | 0.9886 | 0.9771 | 0.9771 |
| FordA | 0.9727 | 0.9272 | 0.9733* | 0.9267 | 0.9267 |
| FordB | 0.9173 | 0.9180 | 0.9186* | 0.9158 | 0.9158 |
| Gun_Point | 1.0000 | 1.0000* | 1.0000* | 1.0000* | 1.0000* |
| Ham | 0.7810 | 0.7714 | 0.8000 | 0.8381* | 0.8000 |
| HandOutlines | 0.9487 | 0.8930 | 0.8870 | 0.9030 | 0.9030 |
| Haptics | 0.5510 | 0.5747* | 0.5584 | 0.5649 | 0.5584 |
| Herring | 0.7030 | 0.7656* | 0.7188 | 0.7500 | 0.7656* |
| InlineSkate | 0.6127 | 0.4655 | 0.5000 | 0.4927 | 0.4927 |
| InsWngSnd | 0.6525 | 0.6616 | 0.6696 | 0.6823* | 0.6818 |
| ItPwDmd | 0.9700 | 0.9631 | 0.9699 | 0.9602 | 0.9708* |
| LrgKitApp | 0.8960 | 0.9200* | 0.9200* | 0.9067 | 0.9120 |
| Lighting2 | 0.8853 | 0.8033 | 0.8197 | 0.7869 | 0.7869 |
| Lighting7 | 0.8630 | 0.8356 | 0.9178* | 0.8219 | 0.9178* |
| Mallat | 0.9800 | 0.9808 | 0.9834 | 0.9838 | 0.9842* |
| Meat | 1.0000 | 0.9167 | 1.0000* | 0.9833 | 1.0000* |
| MedicalImages | 0.7920 | 0.8013 | 0.8066* | 0.7961 | 0.7961 |
| MidPhxAgeGp | 0.8144 | 0.8125 | 0.8150 | 0.8175* | 0.8075 |
| MidPhxCorr | 0.8076 | 0.8217 | 0.8333 | 0.8400 | 0.8433* |
| MidPhxTW | 0.6120 | 0.6165 | 0.6466 | 0.6466* | 0.6316 |
| MoteStrain | 0.9500 | 0.9393 | 0.9569* | 0.9361 | 0.9361 |
| NonInv_Thor1 | 0.9610 | 0.9654 | 0.9657 | 0.9751 | 0.9756* |
| NonInv_Thor2 | 0.9550 | 0.9623 | 0.9613 | 0.9664 | 0.9674* |
| OliveOil | 0.9333 | 0.8667 | 0.9333 | 0.9333 | 0.9667* |
| OSULeaf | 0.9880 | 0.9959* | 0.9959* | 0.9959* | 0.9917 |
| PhalCorr | 0.8300 | 0.8368 | 0.8392* | 0.8380 | 0.8357 |
| Phoneme | 0.3492 | 0.3776* | 0.3602 | 0.3671 | 0.3623 |
| Plane | 1.0000 | 1.0000* | 1.0000* | 1.0000* | 1.0000* |
| ProxPhxAgeGp | 0.8832 | 0.8927* | 0.8878 | 0.8878 | 0.8927* |
| ProxPhxCorr | 0.9180 | 0.9450* | 0.9313 | 0.9313 | 0.9381 |
| ProxPhxTW | 0.8150 | 0.8350 | 0.8275 | 0.8375* | 0.8375* |
| RefDev | 0.5813 | 0.5813 | 0.5947* | 0.5840 | 0.5840 |
| ScreenType | 0.7070 | 0.6693 | 0.7073 | 0.6907 | 0.6907 |
| ShapeletSim | 1.0000 | 0.9722 | 1.0000* | 0.9833 | 0.9833 |
| ShapesAll | 0.9183 | 0.9017 | 0.9150 | 0.9183 | 0.9217* |
| SmlKitApp | 0.8030 | 0.8080 | 0.8133* | 0.7947 | 0.8133* |
| SonyAIBOI | 0.9850 | 0.9817 | 0.9967 | 0.9700 | 0.9983* |
| SonyAIBOII | 0.9620 | 0.9780 | 0.9822* | 0.9748 | 0.9790 |
| StarlightCurves | 0.9796 | 0.9756 | 0.9763 | 0.9767 | 0.9767 |
| Strawberry | 0.9760 | 0.9838 | 0.9864 | 0.9838 | 0.9865* |
| SwedishLeaf | 0.9664 | 0.9792 | 0.9840 | 0.9856* | 0.9856* |
| Symbols | 0.9668 | 0.9839 | 0.9849 | 0.9869 | 0.9889* |
| Synth_Cntr | 1.0000 | 0.9933 | 1.0000* | 0.9900 | 0.9900 |
| ToeSeg1 | 0.9737 | 0.9825 | 0.9912* | 0.9868 | 0.9868 |
| ToeSeg2 | 0.9615 | 0.9308 | 0.9462 | 0.9308 | 0.9308 |
| Trace | 1.0000 | 1.0000* | 1.0000* | 1.0000* | 1.0000* |
| Two_Patterns | 1.0000 | 0.9968 | 0.9973 | 0.9968 | 0.9968 |
| TwoLeadECG | 1.0000 | 0.9991 | 1.0000* | 0.9991 | 1.0000* |
| uWavGest_X | 0.8308 | 0.8490 | 0.8498 | 0.8481 | 0.8504* |
| uWavGest_Y | 0.7585 | 0.7672* | 0.7661 | 0.7658 | 0.7644 |
| uWavGest_Z | 0.7725 | 0.7973 | 0.7993 | 0.7982 | 0.8007* |
| uWavGestAll | 0.9685 | 0.9618 | 0.9609 | 0.9626 | 0.9626 |
| Wafer | 1.0000 | 0.9992 | 1.0000* | 0.9981 | 0.9981 |
| Wine | 0.8890 | 0.8704 | 0.8890 | 0.9074* | 0.9074* |
| WordsSynonyms | 0.7790 | 0.6708 | 0.6991 | 0.6677 | 0.6677 |
| Worms | 0.8052 | 0.6685 | 0.6851 | 0.6575 | 0.6575 |
| WormsTwoClass | 0.8312 | 0.7956 | 0.8066 | 0.8011 | 0.8011 |
| yoga | 0.9183 | 0.9177 | 0.9163 | 0.9190 | 0.9237* |
| Count | - | 43 | 65 | 51 | 57 |
| MPCE | - | 0.0318 | 0.0283 | 0.0301 | 0.0294 |
| Arith. Mean | - | - | 2.1529 | - | 2.5647 |
| Geom. Mean | - | - | 1.8046 | - | 1.8506 |

TABLE II: Wilcoxon Signed Rank Test comparison of each Model

| | WEASEL | 1-NN DTW CV | 1-NN DTW | BOSS | Learning Shapelet | TSBF | ST | EE | COTE | MLP | CNN | ResNet | LSTM-FCN | F-t LSTM-FCN | ALSTM-FCN |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| WEASEL | | | | | | | | | | | | | | | |
| 1-NN DTW CV | 2.39E-10 | | | | | | | | | | | | | | |
| 1-NN DTW | 2.53E-12 | 7.20E-04 | | | | | | | | | | | | | |
| BOSS | 4.27E-03 | 1.82E-07 | 5.31E-11 | | | | | | | | | | | | |
| Learning Shapelet | 2.00E-04 | 2.53E-02 | 2.33E-04 | 1.94E-02 | | | | | | | | | | | |
| TSBF | 2.18E-05 | 1.59E-01 | 2.49E-03 | 4.36E-03 | 4.73E-01 | | | | | | | | | | |
| ST | 1.29E-01 | 1.05E-07 | 9.64E-11 | 2.39E-01 | 1.61E-03 | 3.60E-04 | | | | | | | | | |
| EE | 4.51E-05 | 3.45E-07 | 1.31E-10 | 1.37E-02 | 6.13E-01 | 2.02E-01 | 1.39E-03 | | | | | | | | |
| COTE | 5.44E-01 | 3.05E-14 | 3.03E-16 | 6.21E-04 | 4.76E-07 | 1.13E-06 | 4.24E-03 | 3.54E-11 | | | | | | | |
| MLP | 2.56E-07 | 5.21E-01 | 3.41E-01 | 6.89E-05 | 1.44E-02 | 8.37E-02 | 6.76E-06 | 4.88E-03 | 2.84E-08 | | | | | | |
| FCN | 2.77E-01 | 1.84E-10 | 2.14E-15 | 1.03E-03 | 3.65E-06 | 1.54E-06 | 8.85E-03 | 6.07E-06 | 4.82E-01 | 2.79E-09 | | | | | |
| ResNet | 5.67E-01 | 1.82E-10 | 5.95E-15 | 4.38E-03 | 1.32E-05 | 3.56E-06 | 2.47E-02 | 1.09E-05 | 9.61E-01 | 4.64E-08 | 2.52E-01 | | | | |
| LSTM-FCN | 4.92E-06 | 1.92E-17 | 8.59E-21 | 3.00E-11 | 2.65E-12 | 4.04E-12 | 9.93E-13 | 5.14E-13 | 1.60E-07 | 1.61E-14 | 1.05E-07 | 4.91E-10 | | | |
| F-t LSTM-FCN | 1.23E-08 | 5.17E-19 | 5.77E-22 | 3.35E-13 | 2.20E-13 | 1.12E-13 | 3.44E-14 | 1.25E-14 | 2.81E-10 | 5.09E-16 | 3.35E-12 | 4.58E-15 | 7.53E-05 | | |
| ALSTM-FCN | 1.34E-07 | 2.74E-18 | 5.14E-21 | 1.34E-12 | 3.38E-12 | 7.48E-13 | 7.11E-14 | 1.26E-13 | 1.30E-08 | 1.70E-15 | 3.74E-09 | 1.33E-11 | 8.53E-04 | 3.06E-02 | |
| F-t ALSTM-FCN | 4.58E-08 | 1.01E-18 | 1.18E-21 | 1.44E-12 | 2.41E-12 | 4.63E-13 | 3.96E-14 | 4.12E-14 | 2.56E-09 | 1.87E-15 | 2.60E-10 | 1.12E-12 | 5.96E-05 | 1.89E-01 | 5.40E-02 |

TABLE III: Summary of advantages of the proposed models

| Advantage | LSTM-FCN | F-t LSTM-FCN | ALSTM-FCN | F-t ALSTM-FCN |
|---|---|---|---|---|
| Performance | | ✓ | | ✓ |
| Visualization | | | ✓ | ✓ |

model outperforms the Phase 1 LSTM-FCN model. One of the major advantage of using the Attention LSTM cell is it provides a visual representation of the attention vector. The Attention LSTM also benefits from *fine-tuning*, but the effect is less significant as compared to the general LSTM model. A summary of the performance of each model type on certain characteristics is provided on Table III.

## V. CONCLUSION & FUTURE WORK

With the proposed models, we achieve a potent improvement in the current state-of-the-art for time series classification using deep neural networks. Our baseline models, with and without fine-tuning, are trainable end-to-end with nominal preprocessing and are able to achieve significantly improved performance. LSTM-FCNs are able to augment FCN models, appreciably increasing their performance with a nominal increase in the number of parameters. ALSTM-FCNs provide one with the ability to visually inspect the decision process of the LSTM RNN and provide a strong baseline on their own. *Fine-tuning* can be applied as a general procedure to a model to further elevate its performance. The strong increase in performance in comparison to the FCN models shows that LSTM RNNs can beneficially supplement the performance of FCN modules for time series classification. An overall analysis of the performance of our model is provided and compared to other techniques.

There is further research to be done on understanding why the attention LSTM cell is unsuccessful in matching the performance of the general LSTM cell on some of the datasets. Furthermore, extension of the proposed models to multivariate time series is elementary, but has not been explored in this work.

## REFERENCES

[1] M. W. Kadous, "Temporal Classification: Extending the Classification Paradigm to Multivariate Time Series," *New South Wales, Australia*, 2002.

[2] J. Lin, E. Keogh, L. Wei, and S. Lonardi, "Experiencing SAX: A Novel Symbolic Representation of Time Series," *Data Mining and Knowledge Discovery*, vol. 15, no. 2, pp. 107–144, apr 2007.

[3] M. G. Baydogan, G. Runger, and E. Tuv, "A Bag-of-Features Framework to Classify Time Series," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 11, pp. 2796–2802, nov 2013.

[4] P. Schäfer, "The BOSS is Concerned with Time Series Classification in the Presence of Noise," *Data Mining and Knowledge Discovery*, vol. 29, no. 6, pp. 1505–1530, sep 2014.

[5] P. Schäfer, "Scalable Time Series Classification," *Data Mining and Knowledge Discovery*, vol. 30, no. 5, pp. 1273–1298, 2016.

[6] P. Schäfer and U. Leser, "Fast and Accurate Time Series Classification with WEASEL," *arXiv preprint arXiv:1701.07681*, 2017.

[7] J. Lines and A. Bagnall, "Time Series Classification with Ensembles of Elastic Distance Measures," *Data Mining and Knowledge Discovery*, vol. 29, no. 3, pp. 565–592, jun 2014.

[8] A. Bagnall, J. Lines, J. Hills, and A. Bostrom, "Time-Series Classification with COTE: The Collective of Transformation-Based Ensembles," *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 9, pp. 2522–2535, 2015.

[9] Z. Cui, W. Chen, and Y. Chen, "Multi-Scale Convolutional Neural Networks for Time Series Classification," *arXiv preprint arXiv:1603.06995*, 2016.

[10] Z. Wang, W. Yan, and T. Oates, "Time Series Classification from Scratch with Deep Neural Networks: A Strong Baseline," in *Neural Networks (IJCNN), 2017 International Joint Conference on*. IEEE, 2017, pp. 1578–1585.

[11] Y. Chen, E. Keogh, B. Hu, N. Begum, A. Bagnall, A. Mueen, and G. Batista, "The UCR Time Series Classification Archive," July 2015, www.cs.ucr.edu/~eamonn/time_series_data/.

[12] C. Lea, R. Vidal, A. Reiter, and G. D. Hager, "Temporal Convolutional Networks: A Unified Approach to Action Segmentation," pp. 47–54, 2016.

[13] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," in *International Conference on Machine Learning*, 2015, pp. 448–456.

[14] L. Trottier, P. Giguère, and B. Chaib-draa, "Parametric Exponential Linear Unit for Deep Convolutional Neural Networks," *arXiv*, pp. 1–16, may 2016. [Online]. Available: http://arxiv.org/abs/1605.09332

[15] R. Pascanu, C. Gulcehre, K. Cho, and Y. Bengio, "How to Construct Deep Recurrent Neural Networks," *arXiv preprint arXiv:1312.6026*, 2013.

[16] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[17] A. Graves *et al.*, *Supervised Sequence Labelling with Recurrent Neural Networks*. Springer, 2012, vol. 385.

[18] D. Bahdanau, K. Cho, and Y. Bengio, "Neural Machine Translation by Jointly Learning to Align and Translate," *arXiv preprint arXiv:1409.0473*, 2014.

[19] M. Lin, Q. Chen, and S. Yan, "Network in Network," *arXiv preprint arXiv:1312.4400*, 2013.

[20] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting." *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.

[21] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How Transferable are Features in Deep Neural Networks?" in *Advances in neural information processing systems*, 2014, pp. 3320–3328.

[22] G. King and L. Zeng, "Logistic Regression in Rare Events Data," *Political analysis*, vol. 9, no. 2, pp. 137–163, 2001.

[23] D. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[24] K. He, X. Zhang, S. Ren, and J. Sun, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on Imagenet Classification," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.