N Sasipriyaa[1], Dr.P.Natesan[2], K.Venu [3], S Mohamed Riyaz[4],B Karthikeyan[5], E Kavin Mukil[6]

Research Article

# Handwritten Ancient Tamil Character Recognition Using Generative AdversarialNetwork (Gan)

**N Sasipriyaa[1], Dr.P.Natesan[2], K.Venu [3], S Mohamed Riyaz[4],B Karthikeyan[5], E Kavin Mukil[6]**

[1*] Assistant Professor(Srg), Department Of Computer Science And Engineering, Kongu Engineering College,
Sasipriya@Kongu.Ac.In
(Corresponding Author)
[2] Professor, Department Of Computer Science And Engineering, Kongu Engineering College
Natesanp@Kongu.Ac.In

[3]assistant Professor, Department Of Computer Science And Engineering, Kongu Engineering College,
Venu.Cse@Kongu.Ac.In
[4,5,6] Ug Scholars, Department Of Computer Science And Engineering, Kongu Engineering College

## Abstract

Recognition Of Handwritten English Characters Has Been Done By Peoples Of Various Parts Of The Geographical Location. Although The Recognition Of These Characters Has Been Achieved For Many Languages Other Than English, It Has Not Yet Been Achieved For Indian Languages Such As Tamil, Kannada, Telugu And Etc. Ancient Period Tamil Characters Are Different From Current Tamil Characters. People Have Shaped Their Discoveries In Fields Such As Medicine, Astronomy, Art And Science. The People Who Have Knowledge In Recognizing Ancient Characters Are Few. It Made The Situation To Recognize Characters That Are Written In Various Sculptures And Temples. We Will Learn About Their Lifestyle By Recognising The Letters. This Model Helps In Recognizing Tamil Characters Used In The Period Of A.D 3rd And 4th Century. The Aim Of The Proposed System Is To Recognize Handwritten Ancient Period Tamil Characters Using Gan Based Convolutional Neural Networks (Cnn). Generative Adversarial Networks (Gans)Are An Exciting Recent Innovation In Deep Learning. Gan Contains Generative Model Which Create New Data Instances That ResembleThe Training Data. Convolution Neural Network Is Deep Learning Models Which Takes Input Generated By Gan And Process It With Its Multiple Layers To Recognize The Characters. The Implementation Of Such Training Model Results With The Accuracy Of 96%.

**Keywords**: Generative Adversarial Networks, Discriminator, Generator, Tamil, Cnn, Deep Learning, Convolutional Neural Networks.

# I Introduction

## 1.1 Character Recognition

Character Recognition Is The Electronic Conversion Of Typed, Handwritten Or Printed Characters From A Scanned Document. The Handwritten Character Recognition Can Be Used In Converting Written Documents To Printed Form, Bank Statements, Business Cards. It Is Used To Identify The Characters Of DifferentLanguages To Identify The Exact Meaning. This System Is Used To Recognize The Tamil Characters Of The Ancient Period Which Is Of Ad $3^{rd}$ Century And Ad $4^{th}$ Century. These Ancient Characters Are Found In Sculptures That Have Been Scattered Around Many Parts Of South India MainlyTamilnadu And Karnataka.

### Tamil Language

Tamil Language Is A Language Being Used Mostly In Tamil Nadu, India. It Behaves As A Root Language For Other Languages Such As Kannada, Telugu. Tamil Has A Total Of 247 Characters In Total In Which 12 Are Vowels, 18 Are Consonants, 1 Is A Special Symbol, Others Are Compound Ones From Combination Of Vowels And Consonants. Ancient Characters Are The Starting Part Of Tamil Language In Which The Current Tamil Letters Are Evolved. The Ancient Period Characters Are Difficult To Recognize, Since Few Amount Of Resources Are Available To Know The Exact Meaning Of The Character And Also The Complex Structure Of The Character. Among The 247 Characters In Tamil There Is No Difference For Some Of The Nedil AndKuril Characters In Ancient Period. These Kinds Of Things Made The Ancient Period Characters Quite Hard.

## 1.2 Generative Adversarial Network

Generative Adversarial Networks, Or Gans For Short, Are An Approach To Generative Modeling Using Deep Learning Methods, Such As Convolutional Neural Networks. Generative Modeling Is An Unsupervised Learning Task In Machine Learning That Involves Automatically Discovering And Learning The Regularities Or Patterns In Input Data In Such A Way That The Model CanBe Used To Generate Or Output New Examples That Plausibly Could Have Been Drawn From The Original Dataset. There Are Two Modules Available In Gan. They Are 1. Generator 2. Discriminator

### Generator

The Generator Model Takes A Fixed-Length Random Vector As Input And Generates A Sample In The Domain. The Vector Is Drawn From Randomly From A Gaussian Distribution, And The Vector Is Used To Seed The Generative Process. After Training, Points In This Multidimensional Vector Space Will Correspond To PointsIn The Problem Domain, Forming A Compressed Representation Of The Data Distribution. This Vector Space Is Referred To As A Latent Space, Or A Vector Space Comprised Of Latent Variables. Latent Variables, Or Hidden Variables, Are Those Variables That Are Important For A Domain But Are Not Directly Observable.

### Discriminator

N Sasipriyaa[1], Dr.P.Natesan[2], K.Venu [3], S Mohamed Riyaz[4],B Karthikeyan[5], E Kavin Mukil[6]

The Discriminator Model Takes An Example From The Domain As Input (Real Or Generated) And Predicts A Binary Class Label Of Real Or Fake (Generated). The Real Example Comes From The Training Dataset. The Generated Examples Are Output By The Generator Model. The Discriminator Is A Normal (And Well Understood) Classification Model. After The Training Process, The Discriminator Model Is Discarded As We Are Interested In The Generator.

## Ii Related Work

Md. Mahbubar Rahman, M. A. H. Akhand, Shahidul Islam, Pintu Chandra Shill. "Bangla Handwritten Character Recognition" 2017 [1] Using Convolutional Neural Network . This Paper Implements Cnn To Extract Features From The Input Image Through Its Convolutional Operations. The Cnn Have A Ability To Perform Augmentation To The Training Data. Then By Forcing The Replication Of Weight Configurations To Obtain Feature Map In The Next Layer. By Reducing Spatial Resolution A Certain Degree Of Shift And Distortion Invariance Is Also Achieved. Using Same Set Of Weights For All Features Results In Decreasing Of Free Parameter Significantly. The Cnn Structure Holds 2 Convolutional Layers With 5 X 5 Receptive Fields And 2 Sub Sampling Layers With 2 X 2 Averaging Areas. Input Layer Contains 784 Nodes For 28 X 28 Pixels Image. This Results In Acceptance Accuracy Rate Of 77.7 % .

Manigandan T, Vidhya V, Dhanalakshmi V, Nirmala B "Tamil Character Recognition From Ancient Epigraphical Inscription Using Ocr And Nlp" 2017 [2]- Recognition Of Ancient Tamil Characters Is One Of The Challenging Task For Epigraphers As The Language Has Evolved With Different Characters Set. This Proposed Work Mainly Focuses On Recognition Of Various Tamil Characters Between 9th And 12th Centuries Using Ocr And Nlp Techniques. In This Work, The Inscription Images Collected From Tamil Nadu, Archaeological Department Are Pre-Processed And Segmented. During The Segmentation Process The Colour Images Were Converted To Gray Image And To Binary Image Based On Threshold Value. From Segmented, Image Features Like Number Of Lines, Curves, Loops And Dots Have Been Extracted Using Scale Invariant Feature Transform (Sift) Algorithms For Each Letter To Identify The Exact Character. Characters Will Be Classified And Constructed Based On Vectors Extracted, Using Support Vector Machine (Svm) Classifier And The Patterns Of The Character Will Be Matched With Known Characters And Predicted Using Trigram Technique. Each Identified Character Will Be Assigned With Its Corresponding Unicode Value And It Will Be Updated In The Image Corpus For Further Character Identification, And To Make The System In Identifying The Characters More Effectively. Thus, The Proposed System Can Solve The Major Problems In Reading The Inscription Images.

R. Jayakanthan1, A. Hiran Kumar2, N. Sankarram3, B. S. Charulatha4, Ashwin Ramesh5 "Handwritten Tamil Character Recognition Using Resnet" 2020 [3]. In This Paper Initially The Data Set Are Different In Sizes Then The Data Augmentation Is Performed To Convert The Irregular Pixel Image Into 128 X 128 Definite Pixel Images. Multistage Character Recognition Scheme Is Used To Avoid System Failure In Small Dataset. Before Training The Image ,Image PreProcessing Is Done To Avoid Unwanted Noises. Then The Images Are Passed Into Resnet Residual Neural Network To Train Hundred And Thousands Of Layer To Get Results In Efficient Manner. The Process Of Transfer Learning And Fine Tuning Are Done To Make The Process To Done Efficiently.

N. Prameela, P. Anjusha, R. Karthik "Off-Line Telugu Handwritten Characters Recognition UsingOptical Character Recognition" 2017 [4]. The Purpose Of The Proposed Paper Is To Detect Telugu Handwritten Offline Characters Using Optical Character Recognition, Ocr Is One Of The Most Popular And Challenging Patterns Of Pattern Recognition. This Paper Proposes An Ocr System Of Telugu Texts Consisting Of Three Phases, Namely Pre-Processing, Feature Extracting, And Classification. In The Advanced Progression Phase, We Applied The Filter Between The Input Characters And Used The Standard Method And Skeletonization Over The Characters To Extract Pixel Points On The Edge Of The Border. In The Feature Release Phase, First Each Character Is Divided Into $3 \times 3$ Grids And The Centroid Corresponding To All Nine Sections Is Tested. With This We Can See Characters Of Different Styles. After That, We Pulled A Horizontal And Vertical Guessing Angel Into A Nearby Pixel Character Called Binary External Symmetry Axis Constellation With An Unrestricted Handwritten Character Where We Calculate The Vertical And Vertical Distance Of The Euclidean By The Same Pixel Approximate From The Center Of Each Location. Then We Calculate The Euclidean Distance And The Angular Values Of The Zones. These Are Considered To Be The Core Values Of Our Proposed System. Finally, Both Vector Support Systems (Svm) And Quadratic Discriminate Classifier (Qda) Were Used Separately As A Separator.

## Iii Proposed System

System For Tamil Handwritten Character Recognition Consists Of A Command Line Utility Which Can Train The Network And Allows The User To Predict TheCharacter In The Image. System Also Consists Of Several Pre-Processing Modules Which Are Used For Pre-Processing Of Training And Testing Data. The User Data Will Be Pre-Processed On The Fly Before It Gets Fed IntoThe System. The System Is Designed To Learn Vowels Of Tamil Language. There Are 12 Vowels In Tamil Language. Initially The Model For Character Recognition Was Done With Convolutional Neural Network. This Model Uses Convolutional Network Along With Adam Optimizer And The Final Result Is Derived From Softmax Function. Three Convolution Layers With Filter Sizes Of 8,16,32 A Max-Booling With The Size Of 2x2 Matrix. The Output Activation Function Has Been Given To Flatten With The Size Of 4096 And The Given To Next Dense Layers. Dropout Is Used To Reduce Overfit Then Finally Neurons Of 26 Classes And Activation Function As Softmax. Since We Have Limited Number Of Images For The Training, The Accuracy Was Limited To 82.62%.

### 3.1 Data Augmentation

Data Augmentation Is The Significant Process It Is Used In Regularization Process To Train The Data In Deep Learning Model. The Data Augmentation Is Done After Resizing The Images Into Same Pixel Level. The Augmentation Provides Similar Images With Respect To Original Images Instead Of Same Image. The Irregular Pixel Images In The System Are Converted Into 64 X 64 Definite Pixel Images.
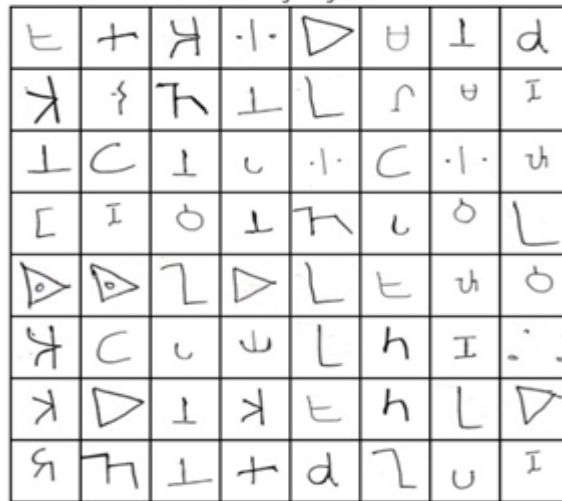
N Sasipriyaa[1], Dr.P Natesan[2], K Venu [3], S Mohamed Riyaz[4] R Karthikeyan[5], E Kavin Mukil[6]

**Fig. 1 Training Images**

## Iv Methodology

The Ancient Period Handwritten Characters Are Given As Input For The System To Train The Model. The Dataset Is Further Processed In Generator Module Of The Gan To Produce Large Amount Of Data. We Are Using Gan To Produce A Greater Number Of Datasets Because Of The Limited Amount Of Resources Available. The Real Images And The Generated Images Are Further Discriminated Using Discriminator And Passed To The Training Module To Train The Dataset. After Training The Data Testing The Data Is Done To Check The Accuracy Of The Model.
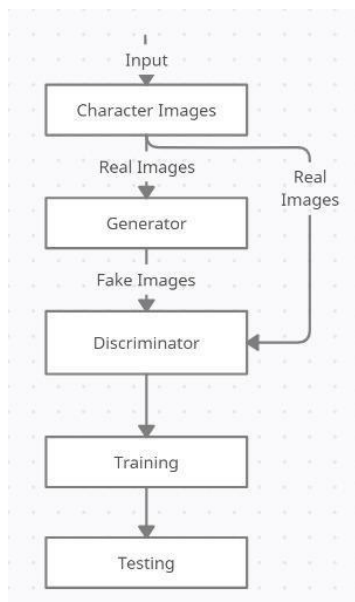


Fig. 2 System Diagram

### 4.1 Image Acquisition

**4.1.1 Data Pre-Processing**

Data Pre-Processing Plays A Major Role In Any Natural Language Processing Task. The Good Pre-Processed Data Yields The Good Result In Any System. Here The Ancient Period Handwritten Tamil Characters Are Considered As Input Data. The Structure Of The Characters May Vary From One Another. Some Characters

Does Not Have Different Meanings But Same Structure. Since The Model Is For Handwritten Character Recognition The InputImage Size And Character Structure May Vary Depends On The Person Using.Hand Written Ancient  Character Are Very Limited In Numbers And Very Few Amounts Of Resources Only Available For Processing It To The System. The Persons Who Know The Exact Meaning And Tend To Translate It Into The Current Tamil Characters Are Very Few In Numbers. The Size Of The Input Images Varies From Person To Person Who Writing The Tamil Ancient Characters. The System Model Input Needs To Be Of Same Size To Process The Inputs So Images Tend To Be Resized AtThe Same Amount Of Parameters.

Since The Resources Available Are Very Less AndWhich Is Not Enough For Training The Module We Are UsingGenerator Module Of Gan To Generate A Greater NumberOf Images By Adding Different Noise Level To The Real Images. Then The Images Generated And The Real Images Are Combined Together To Train The Model.
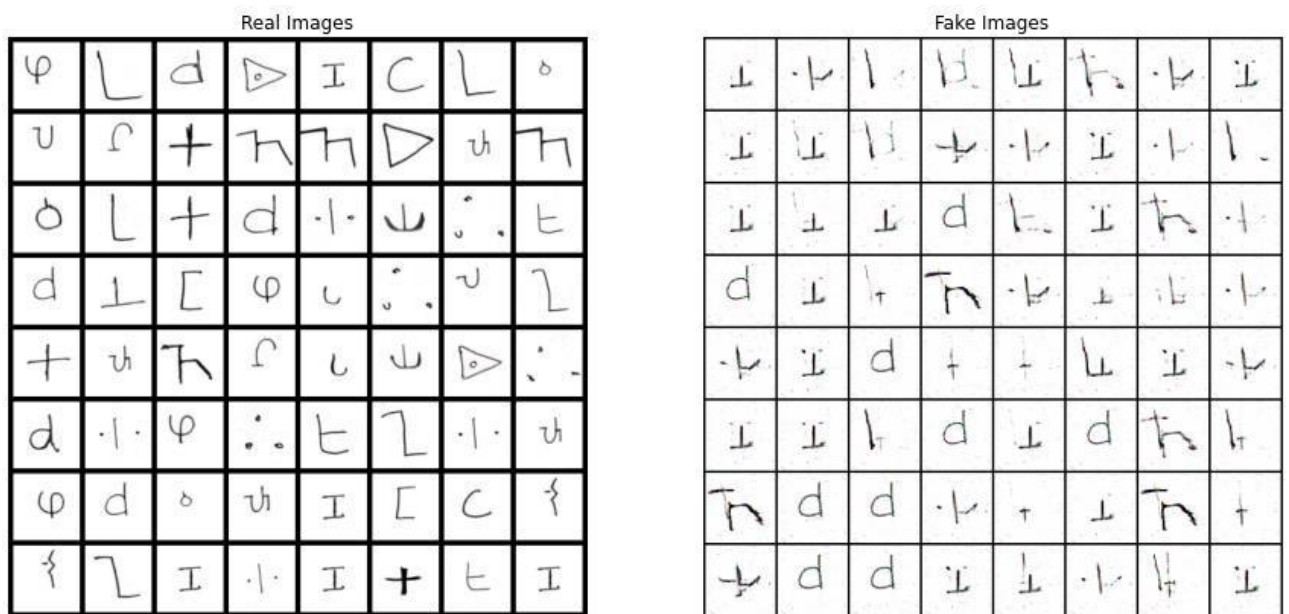


**Fig 3. Real Images And Generated Images**

**4.1.2. Training Dataset**

Generative Adversarial Networks Is One Of The Main Key Features In Hand Written Ancient Character Recognition. This Makes The System Possible To Train Hundreds And Thousands Of Images From Small Number Of Images As Resources. Generally, Most Of The Neural Networks Train The Model Using The Back Propagation. When A Smaller Number Of Datasets Are Used To Train The Model There Is A Chance Of Failure State To Attain. Since For The Handwritten Tamil Ancient Character Recognition
Only A Limited Number Of Resources Are Available We Are Preparing The Dataset Using Generator Module Of The Gan Which Helps In Producing Large Number Of Images From The Limited Number Of Real Images Prepared Manually.

In Our Approach, We Have Used Generative Adversarial Network In Which The Generator Module Helps In Producing The Larger Amount Of Dataset From TheLimited Available Resources. Which In-Turn Will Help To Train The Data With Larger Amount Of Dataset Produced Using Generator Model In Gan.

N Sasipriyaa[1], Dr.P.Natesan[2], K.Venu [3], S Mohamed Riyaz[4],B Karthikeyan[5], E Kavin Mukil[6]

After Generating The Dataset Through The Generator, It Has Been Used To Train The Data For The Model. The Classification Is Done In The Discriminator Module Of The Gan To Assign The Data Produced As Per The Classes. Conv2d Layer Is Used For Classifying The Given Input In The Discriminator. This Layer Creates A Convolutional Kernel That Is Convolved With The Layer Input To Produce A Tensor Of Outputs. If Use-Bias Is True,A Bias Vector Is Created And Added To The Outputs. Finally,If Activation Is Not None, It Is Applied To The Outputs As Well.

The Leaky Relu (Lrelu Or Lrel) Modifies The Function To Allow Small Negative Values When The Input Is Less Than Zero. The Leaky Rectifier Allows For A Small, Non-Zero Gradient When The Unit Is Saturated And Not Active. In This Way, The Acquisition Layer Acquires Both Past And Future Provinces. Future Data Is Therefore Available From The Current Situation.
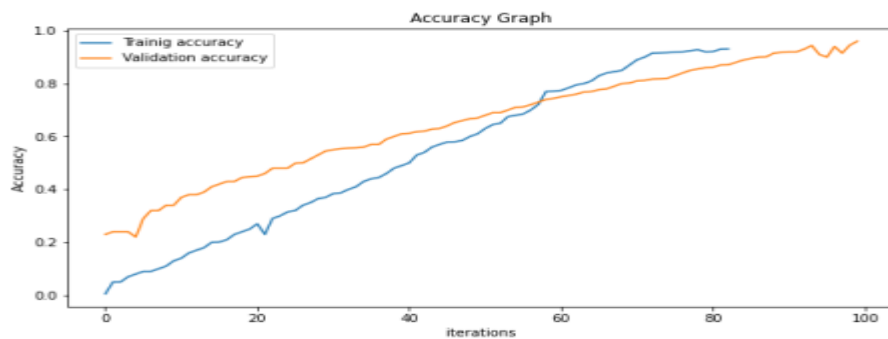


**Fig 4. Generator And Discriminator Loss**

**V Result**

Gan Can Have Two Job Losses: One For Generator Training And One For Racist Training. The Two Losses Work Together To Show The Distance Between The Provision Of Opportunities. In The Loss Schemes, The Generator And The Discriminatory Losses Are Taken At The Same Distance Point Between The Distribution Of Opportunities. In Both Systems, However, The Generator Can Only Affect One Word At A Distance Scale: A Term Indicating The Distribution Of False Data. So During Generator Training We Leave Out Another Name, Which Indicates The Distribution Of Real Data. Generator Loss And Discomfort Look Different In The End, Even If It Is Available In A Single Formula.
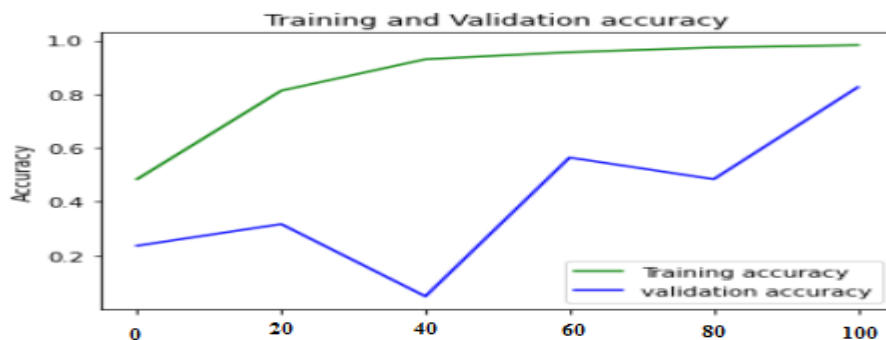


**Fig 5. Accuracy Graph Of Cnn**

As Shown In The Above Accuracy Graphs, Gan Gives Higher Accuracy. Our Proposed Tamil Handwriting Recognition System Is 96.15% Accurate Whereas Recognizing Character Using A Normal Convolution Neural Network Gives The Accuracy Of 86%. This Loss Is Due To A Smaller Number Of Datasets Taken For The Training And Testing. Generator In Gan Model Which Generates Duplicate Images, Enhances The Training And Gives 96.15% Accuracy.

## Vi Conclusion And Future Work

Our Proposed System For Handwritten Tamil Ancient Character Recognition Using Generative Adversarial Network It Includes Dataset And Algorithm. The Dataset Contains More Than 2500 Images That Are Written By Different Persons And Each Image Have Been Resized To 64 X 64 In Dimensions. The Training Images Has Many Characters That Are Visually Similar Which Is Written By Most Of The People. Our Proposed System For Tamil Ancient Character Recognition Gives Accuracy Of 96%. We Have Done Our Proposed System To Check Up To **Uyir** And **Mei_Eluthukal** (30) We Further Extend Our Model To **Uyirmei Eluthukal (216).** The Works Can Be Extended To Recognize Words In Tamil. Thus, Its Become The Repository Of Ancient Period Tamil Character Recognition.

## References

[1]   Md. Mahbubar Rahman, M. A. H. Akhand, Shahidul Islam, Pintu Chandra Shill. "Bangla Handwritten Character Recognition" 2017.

[2]   Manigandan T, Vidhya V, Dhanalakshmi V, Nirmala B "Tamil Character Recognition From Ancient Epigraphical Inscription Using Ocr AndNlp" 2017 .

[3]   R. Jayakanthan1, A. Hiran Kumar2, N. Sankarram3, B. S. Charulatha4, Ashwin Ramesh5"Handwritten Tamil Character Recognition UsingResnet" 2020.

[4]   N. Prameela, P. Anjusha, R. Karthik "Off-Line Telugu Handwritten Characters Recognition Using Optical Character Recognition" 2017.

[5]   Zhengwei Wang, Qi She, Tomas E. Ward "Generative Adversarial Network In Computer Vision: A Survey And Taxonomy" 2020.

[6]   Deepak Dilipkumar, Advisor: Barnab´As P´Oczos "Generative Adversarial Image Refinement For Handwriting Recognition" 2017.

[7]   Md. Mahbubar Rahman, M. A. H. Akhand, Shahidul Islam, Pintu Chandra Shill Dept. Of Computer Science And Engineering Khulna University Of Engineering & Technology Khulna, Bangladesh "Bangla Handwritten Character Recognition Using Convolutional Neural Network"2015.

[8]   Ahmed El-Sawy, Mohamed Loey, Hazem El- Bakry "Arabic Handwritten Characters Recognition Using Convolutional Neural Network"2017

[9]   Mujtaba Husnain 1, Malik Muhammad Saad Missen 1, Shahzad Mumtaz 1, Muhammad Zeeshan Jhanidr 1 "Recognition Of UrduHandwritten Characters Using Convolutional Neural Network".2019

[10]  N. Prameela, P. Anjusha, And R. Karthik "Off-LineTelugu Handwrittencharacters Recognition Using Optical Character Recognition" Ieeetransactions On International Conference On Electronics, Communicationand

N Sasipriyaa[1], Dr.P.Natesan[2], K.Venu [3], S Mohamed Riyaz[4],B Karthikeyan[5], E Kavin Mukil[6]

Aerospace Technology, Iceca 2017 .

[11] M. A. Pragathi, K. Priyadarshini, S. Saveetha, A. Shavar Banu, K. O.Mohammed Aarif"Handwritten Tamil Character Recognition Usingdeep Learning" Ieee Transactions On International Conference Onvision Towards Emerging Trends In Communication And Networking(Vitecon), 2019.

[12] D. Afchar, V. Nozick, J. Yamagishi, And I. Echizen, "Mesonet: A Compact Facial VideoForgery Detection Network," In 2018 Ieee International Workshop On Information Forensics And Security (Wifs). Ieee, 2018.