

# Adaptation of Manipuri OCR FROM Bengali TESSERACT – OCR

---

Summer Internship 2019  
IIT Guwahati  
Dept. Computer Science and Engineering



Sanajaoba Thongram  
NIT Manipur, CSE  
Ningthoujam Justwant Singh  
NIT Manipur, CSE

Under the Guidance of:

Sanasam Ranbir Singh  
Associate Professor

Jennil Thiyam  
PhD Research Scholar

Department of Computer Science and Engineering  
Indian Institute of Technology Guwahati

# ACKNOWLEDGEMENT

---

It is with a sense of gratitude, we acknowledge the efforts of entire hosts of well-wishers who have in some way or other contributed in their own special ways to the success and completion of this summer internship project.

First of all, we express our sage sense of gratitude and indebtedness to our internship supervisor and project guide **Assoc. Prof. Sanasam Ranbir Singh** of **IIT-Guwahati, Department of Computer Science and Engineering**, for his immense support and faith upon us.

We would also like to thank our project mentor **Mr. Jennil Thiyam, Ph.D. Scholar, IIT Guwahati, Centre For Linguistic Science and Technology**, for his kind advice, suggestions and constant help in a lot of ways during the course of this summer internship, without which we wouldn't be able to complete this project.

Further, we would like to express our gratitude to all the scholars and staffs of **OSINT** lab and **CLST, IIT Guwahati** for their valuable suggestions, corporation, motivation, encouragement, and support in many ways which help in the completion of this summer internship project.

# Contents

---

1. Introduction .....	1
2. Bengali Script and Manipuri .....	2
2.1 Vowels, Consonants and Numerals .....	2-4
2.2 Conjuncts in Manipuri .....	5
3. Tesseract – OCR	
3.1 Why Tesseract? .....	6
3.2 Installing Tesseract .....	6
3.3 Fine Tune .....	6
4. Requirements	
4.1 Tools .....	7
4.2 Dependencies .....	7
4.3 Data .....	7
5. Creating Starter Traineddata	
5.1 Modification of ben.training_txt .....	8
5.2 Running tesstrain.sh .....	8
6. Fine Tuning on the new training data .....	9-10
7. Results and Analysis .....	11-12
8. Conclusion and Future work .....	13

# 1. Introduction

---

The Meitei or Manipuri language is one of the scheduled languages (included in the 8th schedule by the 71st amendment of the constitution in 1992). It is a Sino-Tibetan language and the predominant language and lingua franca in the southeastern Himalayan state of Manipur, in northeastern India. It is the official language in government offices. Meitei is also spoken in the Northeast Indian states of Assam and Tripura, and in Bangladesh and Burma (now Myanmar). It is currently classified as a vulnerable language by UNESCO.

Meitei has its own script, which was used until the 18th century. Its earliest use is not known. Pamheiba, the ruler of the Manipur Kingdom who introduced Hinduism, banned the use of the Meitei script and adopted the Bengali script. Between 1709 and the middle of the 20th century, the Meitei language was written using the Bengali script. So, most of the literature and historical records were written in Bengali script. Now in schools and colleges, the Bengali script is gradually being replaced by the Meitei script. And a major transliteration is in need to reproduce those documents in Meitei Mayek script.

There are few [OCR](#) model for Bengali script, unfortunately, none (to be specific a Tesseract-OCR model) works for Manipuri (Bengali script) with acceptable accuracy. An efficient OCR-model could be helpful in speeding up the transliteration and archival process.

And to develop an efficient OCR system for Manipuri (Bengali script) is the objective of our project. Due to similarity between the script used for Bengali and Manipuri, the OCR model (Manipuri) will be adapted on the existing Bengali – Tesseract – OCR model using fine tuning.

## 2. Bengali Script and Manipuri

---

### Vowels

In Bengali Alphabet, there is a total of 11 vowel symbols, whereas in the Bengali – Manipuri alphabet there are 14 vowel symbols in total.

Manipuri Vowels:

অ	আ	ই	ঈ
উ	ঊ	ঋ	৐
এ	ঐ	ও	ঔ
	অং	অঃ	

Bengali Vowels:

অ	আ	ই	ঈ
উ	ঊ	ঋ	
এ	ঐ	ও	ঔ

## Consonants

In Bengali Alphabet, there is a total of 35 consonant symbols, whereas in the Bengali Manipuri alphabet there are 41 consonant symbols in total.

Bengali Vowel symbols used in written Manipuri are as follow:

ক	খ	গ	ঘ	ঙ
চ	ছ	জ	ঝ	ঞ
ট	ঠ	ড	ঢ	ণ
ত	থ	দ	ধ	ন
প	ফ	ব	ভ	ম
য	র	ল	ৱ	শ
ঐ	স	হ	ক্ষ	ড়
ঢ়	য়	ং	ঃ	ৎ
		ৎ		

Bengali script Consonants are as follows:

ক	খ	গ	ঘ	ঙ
চ	ছ	জ	ঝ	ঞ
ট	ঠ	ড	ঢ	ণ
ত	থ	দ	ধ	ন
প	ফ	ব	ভ	ম
য	র	ল	শ	ষ
স	হ	ড়	ঢ়	য়

## Numerals

০	১	২	৩	৪
০	১	২	৩	৪
৫	৬	৭	৮	৯
৫	৬	৭	৮	৯
৫	৬	৭	৮	৯

Consonant along with the **diacritic** form of the vowels অ, আ, ই, ঈ, উ, ঊ, ঋ, এ, ঐ, ও and ঔ:

	ô	a	i	ī	u	ū	ṛ	e	oi	o	ou
kô	ক	কা	কি	কী	কু	কূ	ক্	কে	কৈ	কো	কৌ
khô	খ	খা	খি	খী	খু	খূ	খ্	খে	খৈ	খো	খৌ
gô	গ	গা	গি	গী	গু	গূ	গ্	গে	গৈ	গো	গৌ
ghô	ঘ	ঘা	ঘি	ঘী	ঘু	ঘূ	ঘ্	ঘে	ঘৈ	ঘো	ঘৌ
ñô	ঙ	ঙা	ঙি	ঙী	ঙু	ঙূ	ঙ্	ঙে	ঙৈ	ঙো	ঙৌ

[...]

hô	হ	হা	হি	হী	হু	হূ	হ্	হে	হৈ	হো	হৌ
yô	য়	য়া	য়ি	য়ী	য়ু	য়ূ	য়্	য়ে	য়ৈ	য়ো	য়ৌ
rô	ড়	ড়া	ড়ি	ড়ী	ড়ু	ড়ূ	ড়্	ড়ে	ড়ৈ	ড়ো	ড়ৌ
rhô	ঢ়	ঢ়া	ঢ়ি	ঢ়ী	ঢ়ু	ঢ়ূ	ঢ়্	ঢ়ে	ঢ়ৈ	ঢ়ো	ঢ়ৌ

## Conjuncts

Bengali is full of conjunct consonants. Here are some of them:

ক্‌ kkô	ক্টে ktô	ক্তে ktô	ক্‌ kbô	ক্মে kmô	ক্রে krô	ক্লি klô	ক্ষে ksô	ক্ষ্মে ksmô
ক্সে ksô	গ্‌ gdhô	গ্নে gnô	গ্‌ gbô	গ্মে gmô	গ্নে glô	গ্নে ghnô	ক্‌ nkô	উক্ষে ñksô
ঙ্থে nthô	ঙ্‌ ngô	ঙ্‌ nghô	ঙ্মে ñmô	চ্‌ cchô	চ্‌ cchbô	চ্‌ cñô	জ্‌ jjô	জ্‌ jjbô
জ্‌ jjhô	জ্‌ jñô	জ্‌ jnbô	ঞে ñcô	জ্‌ ñchô	জ্‌ ñjhô	ট্‌ ttô	ট্‌ tbô	ণ্‌ ntô
ঠ্‌ nthô	ণ্‌ ndô	ণ্‌ nñô	ণ্মে ñmô	ত্তে ttô	ত্‌ ttbô	ত্থে tthô	ত্নে tnô	ত্‌ tbô
ত্মে tmô	ত্রে trô	দে ddô	দ্‌ ddhô	দে dbô	দ্‌ dbhrô	ন্‌ ntô	ন্‌ ndô	ন্তে ntô
ত্‌ ntbô	ত্নে ntrô	ন্দে ndô	দ্‌ ndhô	ন্‌ nnô	দে nbô	দে nsô	প্‌ ptô	প্‌ ptô
প্‌ pnô	প্‌ ppô	প্‌ plô	প্‌ psô	ফ্‌ phlô	ভে bhrô	ভ্‌ bhlô	ম্‌ mnô	ম্‌ mphô
ম্‌ mbô	ম্‌ mlô	ল্‌ ltô	ল্‌ ldô	ল্‌ lbô	ল্‌ llô	শ্‌ shchô	শ্‌ skô	ষ্টে ştô
ষে şnô	স্‌ skrô	স্‌ stô	স্‌ strô	স্‌ sbô	হে hnô	হ্মে hmô	হ্‌ hbô	হ্‌ hlô

Some of the Manipuri Conjuncts are as follows:

ক্য	ক্র	ক্ল	ক্‌	হ	ক্স	ক্‌	ক	ক্‌
খ্য	খ্র	খ্ল	খ্‌	ফ	জ্য	খ্‌	ক	ক্‌
জ্য	জ্র	জ্ল	জ্‌	রু	অ	খ্‌	ক	ক্‌
ন্য	এ	ন্ল	দ	ত	অ	খ্‌	ক	ক্‌
প্য	প্র	প্ল	শ	ম	অ	খ্‌	ক	ক্‌
ব্য	ব্র	ব্ল	ষ	প	অ	খ্‌	ক	ক্‌
শ্য	শ্র	শ্ল	উ	ম	অ	খ্‌	ক	ক্‌
হ্য	হ্র	হ্ল	ঠ	হ	অ	খ্‌	ক	ক্‌
ঙ্য	ঙ্র	ঙ্ল	ণ	ন	অ	খ্‌	ক	ক্‌
ত্যা	ত্র	ত্ল	ন্	শ	অ	খ্‌	ক	ক্‌
এ	এ	এ	প	ল	অ	খ্‌	ক	ক্‌
ম্‌	ম্‌	ম্‌	প	শ	অ	খ্‌	ক	ক্‌
ষ	ষ	ষ	প	ম	অ	খ্‌	ক	ক্‌



## 3. Tesseract – OCR

---

Tesseract was originally developed at Hewlett-Packard Laboratories Bristol and at Hewlett-Packard Co, Greeley Colorado between 1985 and 1994, with some more changes made in 1996 to port to Windows, and some C++izing in 1998. In 2005 Tesseract was open sourced by HP and the University of Nevada, Las Vegas (UNLV)[1]. Since 2006 it is developed by Google.

The latest (LSTM based) stable version is [4.0.0](#), released on October 29, 2018.

### 3.1 Why Tesseract – OCR?

Tesseract is an OCR engine with support for **unicode (UTF-8)** and the ability to recognize **more than 100 languages** out of the box. It can be trained to recognize other languages.

Tesseract supports **various output formats**: plain text, hOCR (HTML), PDF, invisible-text-only PDF, TSV. The master branch also has experimental support for ALTO (XML) output. And this OCR engine has a higher accuracy than its counterparts.

### 3.2 Installing Tesseract [linux]

To install **Tesseract 4.x** you can simply run the following command on your **Ubuntu 18.xx bionic**:

```
sudo apt install tesseract-ocr
```

To install the **Developer Tools** which can be used for training, run the following command:

```
sudo apt install libtesseract-dev
```

### 3.3 Fine Tune

Fine tuning is one of the options provided by Tesseract engine to start with an existing trained language, train on our own specific additional data. This may work for problems that are close to the existing training data, but different in some subtle way, like a particularly unusual font. May work with even a small amount of training data.

There are two main types of Fine Tuning in Tesseract 4.xx.

1. Fine tuning for specific fonts.
2. Fine tuning for additional symbols.

## 4. Requirements

---

### 4.1 Tools

1. Tesseract 4.1.0-rc4  
<https://github.com/tesseract-ocr/tesseract/releases/tag/4.1.0-rc4>
2. GPL Ghostscript 9.26 (2018-11-20)
3. OCR-evaluation-Tools by UNLV [1].  
<https://github.com/Sanaj2060/ocr-evaluation-tools>

### 4.2 Dependencies

1. A compiler for C and C++: GCC or Clang
2. GNU Autotools: autoconf, automake, libtool
3. pkg-config
4. Leptonica
5. libpng, libjpeg, libtiff

### 4.3 Data

1. tesseract-ocr/tessdata\_best/ben.traineddata  
[https://github.com/tesseract-ocr/tessdata\\_best/blob/master/ben.traineddata](https://github.com/tesseract-ocr/tessdata_best/blob/master/ben.traineddata)
2. tesseract-ocr/tessdata\_best/script/Bengali.traineddata  
[https://github.com/tesseract-ocr/tessdata\\_best/blob/master/script/Bengali.traineddata](https://github.com/tesseract-ocr/tessdata_best/blob/master/script/Bengali.traineddata)
3. langdata\_master  
<https://github.com/tesseract-ocr/langdata.git>
4. langdata\_lstm\_master  
[https://github.com/tesseract-ocr/langdata\\_lstm.git](https://github.com/tesseract-ocr/langdata_lstm.git)
5. Manipuri text from [Naharol gi Thoudang](#), a Manipuri newspaper.  
<https://github.com/Sanaj2060/tesseract-ocr-Manipuri-Bengali/blob/master/data/manipuri.txt>
6. Manipuri Conjunctions from **Maapi Lyrik**, a Manipuri Primer.  
<https://github.com/Sanaj2060/tesseract-ocr-Manipuri-Bengali/blob/master/data/conjunctions.txt>

## 5. Creating Starter Traineddata

A **starter traineddata** file is given during training, and has to be setup in advance. It can contain:

- Config file providing control parameters.
- **Unicharset** defining the character set.
- **Unicharcompress**, aka the recoder, which maps the unicharset further to the codes actually used by the neural network recognizer.
- Punctuation pattern dawg, with patterns of punctuation allowed around words.
- Word dawg. The system word-list language model.
- Number dawg, with patterns of numbers that are allowed.

Bold elements **must** be provided. Others are optional, but if any of the dawgs are provided, the punctuation dawg must also be provided. A new tool: *combine\_lang\_model* is provided to make a **starter traineddata** from a *unicharset* and optional *wordlists*.

### 5.1 Modification of ben.training\_txt

We select few words from the collection of Manipuri text [4.3.5] to make sure that the text to be added has every unique symbols of Manipuri (Bengali Script) and all the conjuncts from the primer [4.3.6], then append it to:

1. langdata\_master/ben/ben.training\_txt [4.3.3]
2. langdata\_lstm\_master/ben/ben.training\_txt [4.3.4]

The config and other optional files are same with the one used to train the **ben.traineddata**.

### 5.2 Running tesstrain.sh

Training data is created using *tesstrain.sh* as follows:

```
src/training/tesstrain.sh --fonts_dir /usr/share/fonts --lang ben
--linedata_only \
--noextract_font_properties --langdata_dir ../langdata \
--fontlist "Bangla Medium" "Lohit Bengali" \
--tessdata_dir ./tessdata --output_dir ~/tesstutorial/ben_man
```

We run *tesstrain.sh* on the modified ben.training\_txt [5.1.1] and [5.1.2], separately using two fonts, Bangla Medium and Lohit Bengali each.

And each output has 6 files and are as follows:

1. ben.Bangla\_Medium.exp0.lstmf
2. ben.Lohit\_Bengali.exp0.lstmf
3. ben.training\_files.txt
4. ben.charset\_size=xyz.txt
5. ben.traineddata
6. ben.unicharset

\* *tesstrain.sh* also generate tiff and box files from the training\_txt.

## 6. Fine Tuning on the new training data

In tesseract-ocr 4.xx we can fine tune a existing model using:

```
[i] training/combine_tessdata -e tessdata/lang.traineddata \
    ~/lang.lstm

[ii] training/lstmtraining --model_output /path/to/output [--max_image_MB
    6000]\
    --continue_from /path/to/existing/model \
    --traineddata /path/to/traineddata/with/new/unicharset \
    --old_traineddata /path/to/existing/traineddata \
    [--perfect_sample_delay 0] [--debug_interval 0] \
    [--max_iterations 0] [--target_error_rate 0.01] \
    --train_listfile /path/to/list/of/filenames.txt
```

[i] and [ii] will give the output as **model\_name.checkpoint**. The flag with its default values:

- `--debug_interval 0`, the trainer outputs a progress report every 100 iterations.
- `--perfect_sample_delay 0`, discards perfect samples if there haven't been that many imperfect ones seen since the last perfect sample. The current default value of zero uses all samples. In practice the value doesn't seem to have a huge effect, and if training is allowed to run long enough, zero produces the best results.
- `--target_error_rate 0.01`, Stop training if the mean percent error rate gets below this value.
- `--max_iterations 0`, Stop training after this many iterations.

**Note:** All the flags are set to its default value for the whole training.

### Combining the output

The lstmtraining program outputs two kinds of checkpoint files:

- `<model_base>_checkpoint` is the latest model file.
- `<model_base><char_error>_<iteration>.checkpoint` is periodically written as the model with the best training error. It is a training dump just like the checkpoint, but is smaller because it doesn't have a backup model to be used if the training runs into divergence.

Either of these files can be converted to a standard traineddata file as follows:

```
training/lstmtraining --stop_training \
    --continue_from ~/tesstutorial/eng_from_chi/base_checkpoint \
    --traineddata ~/tesstutorial/engtrain/eng/eng.traineddata \
    --model_output ~/tesstutorial/eng_from_chi/eng.traineddata
```

This will extract the recognition model from the training dump, and insert it into the `--traineddata` argument, along with the unicharset, recoder, and any dawgs that were provided during training.

We perform 3 different fine tuning operations using different modified training data and existing models:

sl.no	ben.training_txt	old_traineddata	new_traineddata	error_rate	iteration
1	<a href="#">[5.1.1]</a>	tessdata_best/ ben.traineddata <a href="#">[4.3.1]</a>	man_0.traineddata <a href="#">[O-1]</a>	0.006	89800
2	<a href="#">[5.1.2]</a>	tessdata_best/ ben.traineddata <a href="#">[4.3.1]</a>	man_1.traineddata <a href="#">[O-2]</a>	0.01	125700
3	<a href="#">[5.1.1]</a>	tessdata_best/ script/ Bengali.traineddata <a href="#">[4.3.2]</a>	man_2.traineddata <a href="#">[O-3]</a>	0.01	79200

## 7. Results and Analysis

The three models **man\_0**, **man\_1**, **man\_2** and the existing model **ben** and **Bengali** were tested with the same set of input data.

```
Tesseract input.tiff output[.txt] -l model_name [--psm 3] [--oem 0]
```

page segmentation mode:

--psm 3, fully automatic page segmentation, but no OSD. (Default)

--psm 6, Assume a single uniform block of text.

OCR Engine mode:

--oem 0, legacy mode.

--oem 1, lstm engine.

We prepared a ground truth (.txt) and a pdf for the same, using Ghostscript with size -r720x720 -g6120x7920 the pdf is converted into a .tiff file.

The result (out.txt) generate by tesseract 4.1.0-rc4 with each of the three models are evaluate using the **ocr-evaluation-tools** by UNLV[1]. All the five models that we used are shown, along with the psm, oem, model used (-l), and the accuracy rate of the respective model in the following table.

Model (-l)	Font	psm	oem	Accuracy rate (%)
ben	Lohit_Bengali	3	1	88.59
ben	Lohit_Bengali	6	1	90.86
ben	Bangla_Medium	3	1	90.81
ben	Bangla_Medium	6	1	91.20
Bengali	Lohit_Bengali	3	1	90.92
Bengali	Lohit_Bengali	6	1	93.10
Bengali	Bangla_Medium	3	1	92.41
Bengali	Bangla_Medium	6	1	92.82
man_0	Lohit_Bengali	3	1	90.63
man_0	Lohit_Bengali	6	1	92.99
man_0	Bangla_Medium	3	1	92.63
man_0	Bangla_Medium	6	1	92.76
man_1	Lohit_Bengali	3	1	91.87
man_1	Lohit_Bengali	6	1	94.21
man_1	Bangla_Medium	3	1	94.25
man_1	Bangla_Medium	6	1	94.57
man_2	Lohit_Bengali	3	1	91.71
man_2	Lohit_Bengali	6	1	94.11
man_2	Bangla_Medium	3	1	93.49
man_2	Bangla_Medium	6	1	93.91

Here is a sample of the evaluation report generated by OCR-evaluation-tool.

#### UNLV-ISRI OCR Accuracy Report Version 5.1

53897 Characters  
2926 Errors  
94.57% Accuracy

0 Reject Characters  
0 Suspect Markers  
0 False Marks  
0.00% Characters Marked  
94.57% Accuracy After Correction

Ins	Subst	Del	Errors	
0	0	0	0	Marked
246	1703	977	2926	Unmarked
246	1703	977	2926	Total

Count	Missed	%Right	
107	107	0.00	Unassigned
7717	477	93.82	ASCII Spacing Characters
465	5	98.92	ASCII Special Symbols
6	1	83.33	ASCII Digits
1	1	0.00	ASCII Uppercase Letters
4	4	0.00	ASCII Lowercase Letters
45585	1348	97.04	Bengali
12	6	50.00	General Punctuation
53897	1949	96.38	Total

Errors	Marked	Correct-Generated
1604	0	{ঝ}-{ঝ}
430	0	{ }-{<\n>}
96	0	{ }-{ }
57	0	{ ঝ }-{<\n>ঝ}
52	0	{আ}-{আ}
48	0	{ো}-{ো}
18	0	{এ}-{এ}
18	0	{ঝর}-{ঝর}
16	0	{ড়}-{ড়}
12	0	{এ}-{এঁ} //[not EOF]

In the report above most of the error count for characters " ঝ" and " ড়" are due to the misclustering of the characters "ঝ" as (য ি), and "ড়" as (ড ি) and not because of misclassifications by the model. This type of errors are easily fixable during post-processing.

Evaluation reports and input data link [here](#).

## Conclusion and Future work

---

As we can see the newly generated model perform more accurately than the existing Bengali models on the Manipuri words with the model “**man\_1**” [--psm 6 and –fontlist Bangla Medium], [input tiff -r300 (300 dpi)] giving the highest accuracy of **94.54%** and with post-processing the accuracy can be increase upto 97.35% (considering the report shown above). Currently, we are working to provide an optimized resolution for the input image, as the accuracy of the OCR model depends on the DPI of the image.