

# Kickoff: R&D Project Selection & Team Formation

Meeting #2

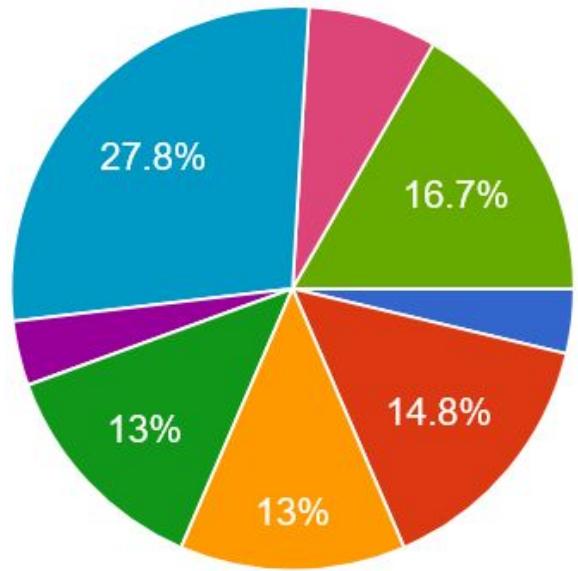
Ye Kyaw Thu  
Lab Leader, LU Lab., Myanmar

15 June 2025 (SUN)

# Overview

1. Survey Results: Top R&D Topics
2. Team Assignments
3. Project Details & Guidance
4. Next Steps & Action Items
5. Open Discussion & Q&A

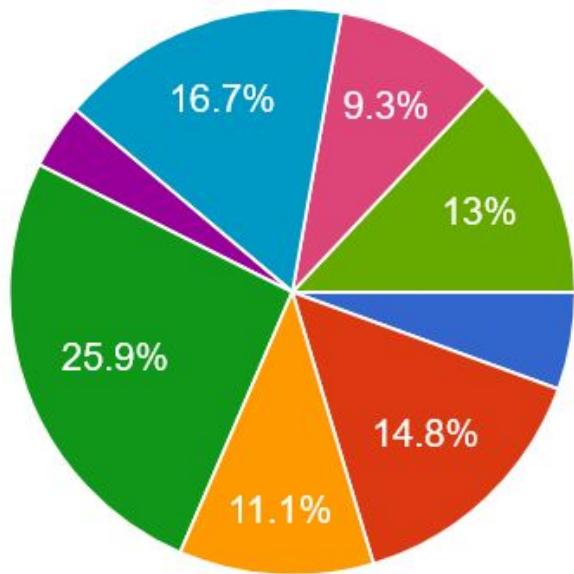
# Survey Results: Top R&D Topics (1st Priority)



- 1. Designing an Educational Programming Language
- 2. Myanmar Text Readability Scoring: Formulas & Evaluation
- 3. Creating a Myanmar SQuAD Dataset...
- 4. Humor Detection in Myanmar Text:...
- 5. Optimizing Myanmar Keyboard Lay...
- 6. Automatic Speech Recognition (AS...)
- 7. Text-to-Speech (TTS) Synthesis for...
- 8. Image Classification or Image2Text

Fig. 1st priority project topics (54 respondents)

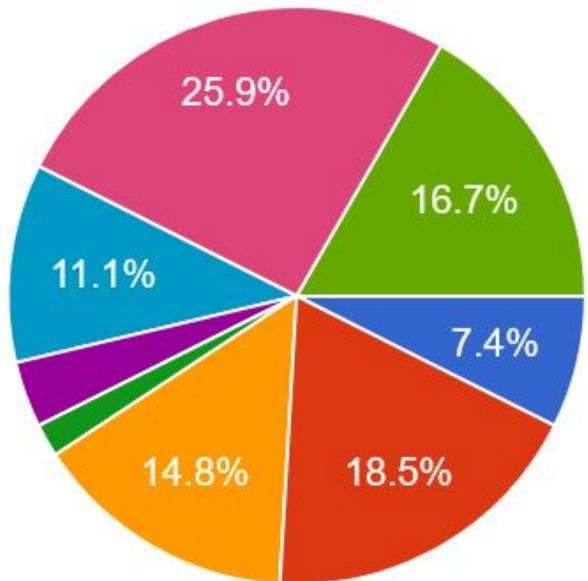
# Survey Results: Top R&D Topics (1st Priority)



- 1. Designing an Educational Programming Language
- 2. Myanmar Text Readability Scoring: Formulas & Evaluation
- 3. Creating a Myanmar SQuAD Dataset...
- 4. Humor Detection in Myanmar Text:...
- 5. Optimizing Myanmar Keyboard Layout...
- 6. Automatic Speech Recognition (ASR)...
- 7. Text-to-Speech (TTS) Synthesis for My...
- 8. Image Classification or Image2Text

Fig. 2nd priority project topics (54 respondents)

# Survey Results: Top R&D Topics (1st Priority)



- 1. Designing an Educational Programming Language
- 2. Myanmar Text Readability Scoring: Formulas & Evaluation
- 3. Creating a Myanmar SQuAD Dataset
- 4. Humor Detection in Myanmar Text:...
- 5. Optimizing Myanmar Keyboard Lay...
- 6. Automatic Speech Recognition (AS...
- 7. Text-to-Speech (TTS) Synthesis for...
- 8. Image Classification or Image2Text

Fig. 3rd priority project topics (54 respondents)

# Survey Results: Top R&D Topics

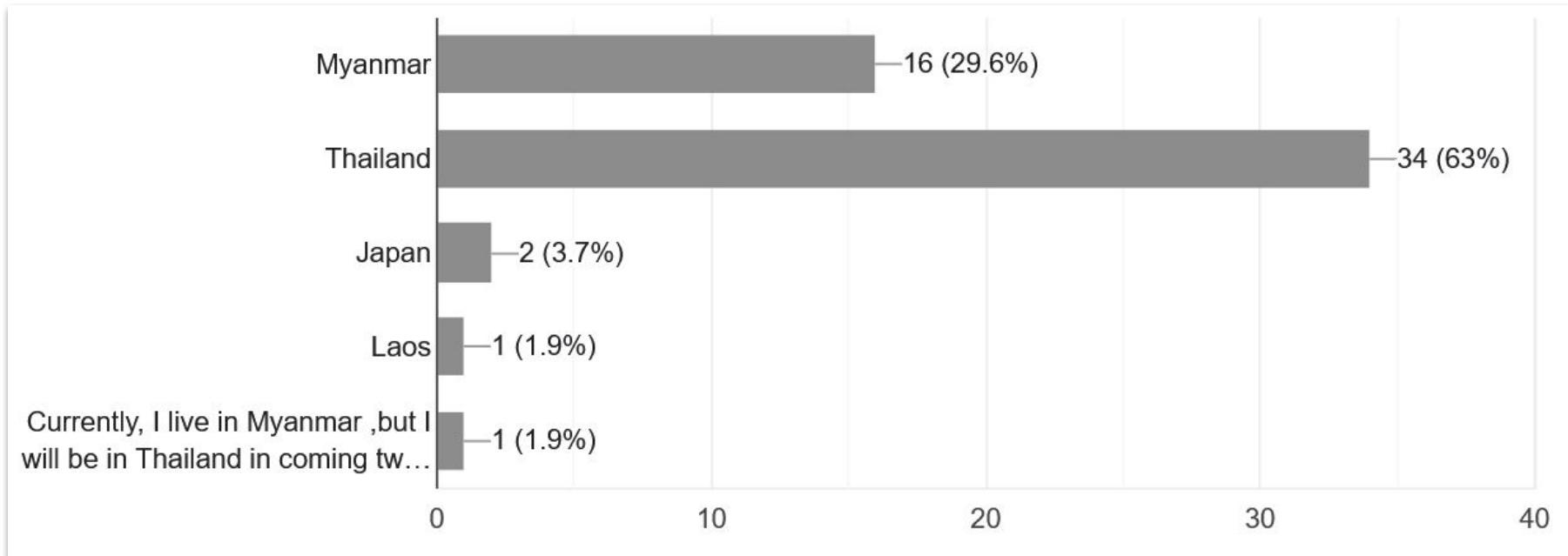


Fig.Current Country of Residence, 54 responses

# Team Assignments (Humor Detection in Myanmar Text)

1. Thura Aung (Thailand)
2. Ni Htwe Aung (Myanmar)
3. Htet Htet Mon (Myanmar)
4. Phyo Thi Khaing (Thailand)
5. Zin Mar Myint (Thailand)
6. Sithu Kyaw Zinn Linn (Thailand)
7. Thant Sin Tun (Myanmar)
8. Kaung Hset Hein (Thailand)

# Team Assignments (ASR)

1. Thura Aung (Thailand)
2. Khaing Hsu Wai (Japan)
3. Sai Wai Yan Phyo (Thailand)
4. Thiri Thaw (Thailand)
5. Moe Chan Myae Maung (Thailand)
6. Thet Htet San (Thailand)
7. Myat Oo Swe (Thailand)
8. Su Sandi Linn (Thailand)
9. Htet Arkar (Thailand)
10. Min Thiha Tun (Myanmar)
11. Saw Zi Dunn (Thailand)
12. Kyi Thant Sin (Thailand)
13. Eaint Lay Hmone (Thailand)
14. Htwe Myat Cho (Thailand)
15. Yadana Myint Hein (Myanmar)

## Team Assignments (TTS)

1. Ye Bhone Lin (Myanmar)
2. Cham Myae Phyo (Thailand)
3. Thant Htut Aung (Myanmar/Thailand)
4. Nann Oak( Caesar) (Thailand)
5. Kyawt Eaindray Win (Thailand)
6. Thet Su Sann (Thailand)
7. Phoo Pwint Cho Thar (Thailand)
8. Htut Ko Ko (Thailand)

## Team Assignments (SQuAD)

1. Hay Man Htun (Lao)
2. Shwe Sin Moe (Myanmar)
3. Khin Mo Mo (Thailand)
4. Chit Ko Ko Lwin (Thailand)
5. Kaung Myat Htet (Thailand)
6. Khaing Zin Theint (Thailand)
7. Thadoe Hein (Thailand)

# Team Assignments (Myanmar Text Readability Scoring)

1. **Khaing Hsu Wai** (Japan)
2. **Hlaing Myat Nwe** (Japan)
3. **Thura Aung** (Thailand)
4. Seng Pan (Thailand)
5. Thiha Nyein (Thailand)
6. Kaung Khant Si Thu (Thailand)
7. Yu Myat Moe (Thailand)
8. Hsu Yee Mon (Thailand)
9. Khin Nyein Chan Thu (Thailand)

# Team Assignments (Image2Text)

1. **Eaint Kay Khaing Kyaw** (Thailand)
2. Than Zaw Toe (Thailand)
3. Pyae Linn (Myanmar)
4. Zayar Htet (Myanmar)
5. Myat Thin Thin Kyi (Myanmar)
6. Thet Hmue Khin (Myanmar)
7. Lynn Myat Bhone (Myanmar)
8. Shin Thant Phyo (Myanmar)
9. Ma Hayman Soe (Myanmar)
10. Khaing Hsu Yee (Myanmar)

# Project Details & Guidance (SQuAD Team)

- SQuAD dataset ကို ရုံးက ဆာဗာပါကို download လုပ်ခဲ့
- Preprocessing တချို့ လုပ်ခဲ့
- Meta ရဲ့ NLLB (No Language Left Behind) NMT နဲ့ အင်္ဂလိပ်ကနေမြန်မာဘာသာကို စပြန်ထား
- ရှုပြုသား NMT ကို သုံးပြုး machine translation လုပ်တာကို အတွေ့အကြီးမရှုတဲ့ သူတွေအနေနဲ့ မြင်သွားစေချင်လို့ practical work တစ်ခုအနေနဲ့ ဆရာလုပ်ပြတာ

# Project Details & Guidance (SQuAD Team)

- HuggingFace ကဲနေ download လုပ်ခဲ့တယ်
- [https://huggingface.co/datasets/rajpurkar/squad/tree/main/plain\\_text](https://huggingface.co/datasets/rajpurkar/squad/tree/main/plain_text)
- Parquet format file နှစ်ဖိုင်

```
(base) ye@1st-gpu-server-197:~/ye/exp/gpt-mt/n11b/data/squad$ file *.parquet
train-00000-of-00001.parquet:      Apache Parquet
validation-00000-of-00001.parquet: Apache Parquet
(base) ye@1st-gpu-server-197:~/ye/exp/gpt-mt/n11b/data/squad$
```

# Project Details & Guidance (SQuAD Team)

- Parquet ဖိုင်က ဒီအတိုင်း ဖွင့်ဖတ်လိုမရဘူး။ အောက်ပါလိုမျိုး  
သိမ်းထားတာမှာ

```
(base) ye@1st-gpu-server-197:~/ye/exp/gpt-mt/n11b/data/squad$ head -n 3 train-00000-of-00001.parquet
PAR1\00000000mL\0000001\5733be284776f41900661182f\07ff\080j\1f\07e\0DF84d058e614000b61bn\B8c\B\bd\fdpc\0d2:\0
8\08F\09j\aj\bj\c\0a640f5\0F\0454J\54F\054J\U\0a70c.\000f64\0MLJ\0n83j85\000bd:p54af\094F\09\0F\05J\U0\ac31.5ff\0F\06\0
94F\00F\07\0T\0d320fe\0F\0F\09\0F\0Y0B\0fY0\0e92101\0F\094F\00F\09PF\09fd3.604\0TF\0n8uJ8\0FT0fb.104uJ\5\0J\5PJ\0J\0u\0b1da66\0
J\0F8u,J\0uJ\51\0b2fe68uhF\0994B89\0F\09PF\0U\0b3d66a\0F\09F\0U0\0b5346pduJ\0uJ\094BTe\0F\095df2\010\0F\0B89F\0YhF\0Y699\0
2\011\08F\0n8YF8YF@\0b7f76\029F\09F\09B\039F\086525c8UJ\0n8u,J8UJ\0U8722 5c9\0F\0a\0F85P\05PF\0a54\087ac.05cb\0p\0
F\0B89\0F\09\08c20cc\08F\0n8uHJ8\0FT8a5cJ\094F894F\00F\09266 \0d85J\0ULF\0\0F\08ULF93160 dJ\0|J\0|J\05J93e6\0dc51J\09P\0
F8U\0J\0U\0J\094c:\000J\09B8tJ\0UhF97460\0e& J\09F89F\09F\0 9816.\00e2\0F\09F\091F\09F\08eb2\0e69F\0YF\094F\094\0a5:\} \0
9\0F\09\0F\09\0F\09b366\0\& J\00TF8nT91c:D\0e\08F\0n89F8a3c6\0\0f3uF\04\0J\0U\05J\05\0a4c52\0\0f\&\0 J\0J\0BT\&\0 \0
5520f3UJ\0uJ\0" J\0XJ\0\0db2(100\0F\09\0F\091F\0b49260c\0J\051F\0dU0J\0\0J\0\082a12\0\0cY\0J\0\0F\0Y\0J\0\083496YJ\0\0F85J\05\0
83e:5UJ\0UJ\0UJ\0UJ\091\0496D\0\0c5J\0J\0<J\0\053:\0\0\0&P\0F\0&P\0J\09BT*\0F\0c006\0\019\0J\0\0J\0c0e6\0\01d51J\0\0B8\0B\0\01a2\0\0
11\0J\0\0F8\0F\0\0292P11b\0TF\0\0TF\0nT\0F8udc:$\011\0\0J\0YL\0F8(J\0(J\0u\0c46\0\011=F\0\0F\0\0F\0D\0\0c742\0\062J\0\0n8F8J\0u\0ca05. \0
62\0H\0J\09F89F\09F\0y\0a:(\02491F\091F\094F\094B\0\0\0cbda6\0\0\0pF8\0pB\0&P\0J\09\0cb2\0127"J\0n8J8yFpd50.\012\0F\0\0B\0\0F\0\0
F\0 06be85543aeaaa\0\08c90MV\0IV\0I\06bf6b\0
\0T6\0
TR8nTiV8i\0\0e46\0 \0V\0-4\0B81\0e86cf60\07n\0-PR88e822T\062V\0n8-P\0bfJ\0T9da7:\094F\0!\0F\0iR\0\0v\0i!\0\08473.,\00d\0
(base) ye@1st-gpu-server-197:~/ye/exp/gpt-mt/n11b/data/squad$
```

# Project Details & Guidance (SQuAD Team)

1. Apache Parquet ဖိုင်ကို normal text file အဖြစ် ဖျေပေးတဲ့ code က ဆရာ GitHub: Tool အောက်မှာ ရှိတယ်။
2. [https://github.com/ye-kyaw-thu/tools/blob/master/python/parquet\\_extractor.py](https://github.com/ye-kyaw-thu/tools/blob/master/python/parquet_extractor.py)
3. \$ python ./parquet\_extractor.py --parquet\_file  
./train-00000-of-00001.parquet
4. \$ python ./parquet\_extractor.py --parquet\_file  
./validation-00000-of-00001.parquet
5. အထက်ပါ command နှစ်ခုနဲ့ ဖြေပြီးသွားရင် original ဖိုင်က CSV ဖိုင်မှိုလို .csv ဖိုင် နှစ်ဖိုင်ကို output အဖြစ် ရလာလိမ်မယ်

# Project Details & Guidance (SQuAD Team)

```
ye@lst-gpu-server-197: ~/ye/exp/gpt-mt/nllb/data/squad
id,title,context,question,answers
5733be284776f41900661182,University_of_Notre_Dame,"Architecturally, the school has a Catholic character. Atop the Main Building's gold dome is a golden statue of the Virgin Mary. Immediately in front of the Main Building and facing it, is a copper statue of Christ with arms upraised with the legend ""Venite Ad Me Omnes"". Next to the Main Building is the Basilica of the Sacred Heart. Immediately behind the basilica is the Grotto, a Marian place of prayer and reflection. It is a replica of the grotto at Lourdes, France where the Virgin Mary reputedly appeared to Saint Bernadette Soubirous in 1858. At the end of the main drive (and in a direct line that connects through 3 statues and the Gold Dome), is a simple, modern stone statue of Mary.",To whom did the Virgin Mary allegedly appear in 1858 in Lourdes France?,"{'text': array(['Saint Bernadette Soubirous']), dtype=object}, 'answer_start': array([515], dtype=int32)}"
5733be284776f4190066117f,University_of_Notre_Dame,"Architecturally, the school has a Catholic character. Atop the Main Building's gold dome is a golden statue of the Virgin Mary. Immediately in front of the Main Building and facing it, is a copper statue of Christ with arms upraised with the legend ""Venite Ad Me Omnes"". Next to the Main Building is the Basilica of the Sacred Heart. Immediately behind the basilica is the Grotto, a Marian place of prayer and reflection. It is a replica of the grotto at Lourdes, France where the Virgin Mary reputedly appeared to Saint Bernadette Soubirous in 1858. At the end of the main drive (and in a direct line that connects through 3 statues and the Gold Dome), is a simple, modern stone statue of Mary.",What is in front of the Notre Dame Main Building?,"{'text': array(['a copper statue of Christ']), dtype=object}, 'answer_start': array([188], dtype=int32)}"
5733be284776f41900661180,University_of_Notre_Dame,"Architecturally, the school has a Catholic character.
```

# Project Details & Guidance (SQuAD Team)

- ဘာသာပြန်ဖို့အတွက်က CSV ဖိုင်ထဲက field တစ်ခုချင်းစီကို ဆွဲထုတ်ယူပို့လိုအပ်တယ်
- Python code ရေးခဲ့တယ်
- Internship ကျောင်းသားတွေလေ့လာလို့ ရွှေအောင် GitHub မှာ repository အသစ်တစ်ခု ဆောက်ပြီး သိမ်းပေးထားမယ်

# Project Details & Guidance (SQuAD Team)

```
#!/usr/bin/env python3
import argparse
import csv
import sys
import os

def clean_text(text):
    """Remove newlines and extra spaces from text"""
    if not text:
        return text
    # Replace newlines with spaces
    text = text.replace('\r', ' ').replace('\n', ' ')
    # Collapse multiple spaces into one
    return ' '.join(text.split())
```

# Project Details & Guidance (SQuAD Team)

```
def extract_field(input_file, column, output_file=None):
    try:
        with open(input_file, 'r', encoding='utf-8') as csvfile:
            # Use csv reader to properly handle quoted fields with commas
            reader = csv.DictReader(csvfile)

            if column not in reader.fieldnames:
                print(f"Error: Column '{column}' not found in CSV file. Available columns: {', '.join(reader.fieldnames)}")
                sys.exit(1)

            data = []
            for row in reader:
                field_value = row[column]
                # Clean the text by removing internal newlines
                cleaned_value = clean_text(field_value)
                data.append(cleaned_value)

        if output_file:
            with open(output_file, 'w', encoding='utf-8') as outfile:
                for item in data:
                    outfile.write(item + '\n')
            print(f"Successfully extracted '{column}' to {output_file}")
        else:
            for item in data:
                print(item)

    except Exception as e:
        print(f"Error processing file: {e}")
        sys.exit(1)
```

# Project Details & Guidance (SQuAD Team)

- clean\_text() မပါဘူး run ရင် အောက်ပါလိုမျိုး ရလဒ်ထွက်လာလိမ့်မယ်။

```
Checking train files... . . .
```

```
...
```

```
(base) · ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ wc · train_*.txt
· 94140 · 714994 · 8706546 · train_answers.txt
· 88128 · 10491130 · 66304332 · train_context.txt
· 87599 · 87599 · 2189975 · train_id.txt
· 87599 · 881343 · 5307142 · train_question.txt
· 87599 · 87599 · 1322786 · train_title.txt
· 445065 · 12262665 · 83830781 · total
...
```

```
...
```

```
original CSV · ဖိုင်နဲ့ နိုင်းယုံကြည့်တဲ့ အခါမှာ . . .
```

```
...
```

```
(base) · ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ wc · train-00000-of-00001.csv
· 94670 · 11915768 · 84323576 · train-00000-of-00001.csv
...
```

# Project Details & Guidance (SQuAD Team)

- SQuAD ဒေတာမှာက အောက်ပါလိုမျိုး field တစ်ခုတည်းမှာပဲ  
တစ်ကြာင်းထက်မကရှိတာမို့လို့

```
20 56beace93aeeaa14008c91e0, Super_Bowl_50, "Super Bowl 50 was an American football game to
determine the champion of the National Football League (NFL) for the 2015 season. The
American Football Conference (AFC) champion Denver Broncos defeated the National Football
Conference (NFC) champion Carolina Panthers 24–10 to earn their third Super Bowl title. The
game was played on February 7, 2016, at Levi's Stadium in the San Francisco Bay Area at
Santa Clara, California. As this was the 50th Super Bowl, the league emphasized the ""golden
anniversary"" with various gold-themed initiatives, as well as temporarily suspending the
tradition of naming each Super Bowl game with Roman numerals (under which the game would
have been known as ""Super Bowl L""), so that the logo could prominently feature the Arabic
numerals 50.", What venue did Super Bowl 50 take place in?", {"text": array(["Levi's
Stadium"]), "Levi's Stadium"}, LF
21      ""Levi's Stadium in the San Francisco Bay Area at Santa Clara""], LF
22      dtype=object), 'answer_start': array([355, 355, 355], dtype=int32)} "LF
23 56beace93aeeaa14008c91e1, Super_Bowl_50, "Super Bowl 50 was an American football game to
determine the champion of the National Football League (NFL) for the 2015 season. The
American Football Conference (AFC) champion Denver Broncos defeated the National Football
```

# Project Details & Guidance (SQuAD Team)

- လက်တွေ field တစ်ခုချင်းစိက္ခ ဆဲထူတ်ဖိုက python code ကို  
တစ်ကြိမ်ထက်မက run ရတာမှုလုံး<sup>ပါ</sup> bash shell script  
(extract\_squad\_fields.sh) ရေးပြုး run တယ်။

```
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/n11b/data/squad$ python ./extract_csv_field.py --help
usage: extract_csv_field.py [-h] --input INPUT --column COLUMN [--output OUTPUT]
```

Extract specific fields from SQuAD CSV files

optional arguments:

```
-h, --help      show this help message and exit
--input INPUT   Input CSV filename
--column COLUMN Column name to extract
--output OUTPUT Output filename (optional)
```

```
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/n11b/data/squad$
```

# Project Details & Guidance (SQuAD Team)

```
#!/bin/bash

# Function to process a single CSV file
process_file() {
    local input_file=$1
    local prefix=$2

    echo "Processing $input_file..."

    # Get all columns from the CSV file
    columns=$(python3 -c "
import csv
with open('$input_file', 'r', encoding='utf-8') as f:
    reader = csv.DictReader(f)
    print(' '.join(reader.fieldnames))
")

    # Extract each column
    for col in $columns; do
        output_file="${prefix}_${col}.txt"
        python3 extract_csv_field.py --input "$input_file" --column "$col" --output "$output_file"
    done
}
```

# Project Details & Guidance (SQuAD Team)

- ပြုးတော့မှ process\_file ဆိုတဲ့ function ကို ခေါ်ပြုး run တဲ့ ပုံစံပါ။

```
# Main script
echo "Starting SQuAD dataset extraction..."

# Process training file
process_file "train-00000-of-00001.csv" "train"

# Process validation file
process_file "validation-00000-of-00001.csv" "valid"

echo "Extraction completed successfully."
```

# Project Details & Guidance (SQuAD Team)

- Extracting all fields with prepared shell script and python code
- \$ time ./extract\_squad\_fields.sh

```
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ wc -l train_*.txt
87599 train_answers.txt
87599 train_context.txt
87599 train_id.txt
87599 train_question.txt
87599 train_title.txt
437995 total
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ wc -l valid_*.txt
10570 valid_answers.txt
10570 valid_context.txt
10570 valid_id.txt
10570 valid_question.txt
10570 valid_title.txt
52850 total
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ _
```

# Project Details & Guidance (SQuAD Team)

- Field ၂ စုစုပေါင်း ၅ခု ရှိတယ် (id,title,context,question,answers)
- Id, title ဆုတ္တံ field နှစ်ခုက မြန်မာလို translation လုပ်စရာ မလိုဘူး

```
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ head -n 3 train_id.txt  
5733be284776f41900661182  
5733be284776f4190066117f  
5733be284776f41900661180  
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ head -n 3 train_title.txt  
University_of_Notre_Dame  
University_of_Notre_Dame  
University_of_Notre_Dame  
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ tail -n 3 valid_id.txt  
5737aaf1c456719005744fd  
5737aaf1c456719005744fe  
5737aaf1c456719005744ff  
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ tail -n 3 valid_title.txt  
Force  
Force  
Force  
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ _
```

# Project Details & Guidance (SQuAD Team)

- မြန်မာလို ဘာသာပြန်ဖို့ အဓိက လိုအပ်တာက context ရယ်။ question ရယ်

```
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/n11b/data/squad$ head -n 3 ./train_context.txt | nl
1 Architecturally, the school has a Catholic character. Atop the Main Building's gold dome is a golden statue of the Virgin Mary. Immediately in front of the Main Building and facing it, is a copper statue of Christ with arms upraised with the legend "Venite Ad Me Omnes". Next to the Main Building is the Basilica of the Sacred Heart. Immediately behind the basilica is the Grotto, a Marian place of prayer and reflection. It is a replica of the grotto at Lourdes, France where the Virgin Mary reputedly appeared to Saint Bernadette Soubirous in 1858. At the end of the main drive (and in a direct line that connects through 3 statues and the Gold Dome), is a simple, modern stone statue of Mary.
2 Architecturally, the school has a Catholic character. Atop the Main Building's gold dome is a golden statue of the Virgin Mary. Immediately in front of the Main Building and facing it, is a copper statue of Christ with arms upraised with the legend "Venite Ad Me Omnes". Next to the Main Building is the Basilica of the Sacred Heart. Immediately behind the basilica is the Grotto, a Marian place of prayer and reflection. It is a replica of the grotto at Lourdes, France where the Virgin Mary reputedly appeared to Saint Bernadette Soubirous in 1858. At the end of the main drive (and in a direct line that connects through 3 statues and the Gold Dome), is a simple, modern stone statue of Mary.
3 Architecturally, the school has a Catholic character. Atop the Main Building's gold dome is a golden statue of the Virgin Mary. Immediately in front of the Main Building and facing it, is a copper statue of Christ with arms upraised with the legend "Venite Ad Me Omnes". Next to the Main Building is the Basilica of the Sacred Heart. Immediately behind the basilica is the Grotto, a Marian place of prayer and reflection. It is a replica of the grotto at Lourdes, France where the Virgin Mary reputedly appeared to Saint Bernadette Soubirous in 1858. At the end of the main drive (and in a direct line that connects through 3 statues and the Gold Dome), is a simple, modern stone statue of Mary.
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/n11b/data/squad$
```

# Project Details & Guidance (SQuAD Team)

- question field ഓല്ലും അരേം

```
(base) ye@1st-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ head -n 15 ./train_question.txt | nl  
1 To whom did the Virgin Mary allegedly appear in 1858 in Lourdes France?  
2 What is in front of the Notre Dame Main Building?  
3 The Basilica of the Sacred heart at Notre Dame is beside to which structure?  
4 What is the Grotto at Notre Dame?  
5 What sits on top of the Main Building at Notre Dame?  
6 When did the Scholastic Magazine of Notre dame begin publishing?  
7 How often is Notre Dame's the Juggler published?  
8 What is the daily student paper at Notre Dame called?  
9 How many student news papers are found at Notre Dame?  
10 In what year did the student paper Common Sense begin publication at Notre Dame?  
11 Where is the headquarters of the Congregation of the Holy Cross?  
12 What is the primary seminary of the Congregation of the Holy Cross?  
13 What is the oldest structure at Notre Dame?  
14 What individuals live at Fatima House at Notre Dame?  
15 Which prize did Frederick Buechner create?  
(base) ye@1st-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ _
```

# Project Details & Guidance (SQuAD Team)

- answer field ခဲ့ format က အောက်ပါအတိုင်းမြှုလို့ မြန်မာစာ အတွက်က သပ်သပ်ပြင်ဆင်ရလို့မြှုမယ်

```
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$ head ./train_answers.txt | nl
 1 {'text': array(['Saint Bernadette Soubirous'], dtype=object), 'answer_start': array([515], dtype=int32)}
 2 {'text': array(['a copper statue of Christ'], dtype=object), 'answer_start': array([188], dtype=int32)}
 3 {'text': array(['the Main Building'], dtype=object), 'answer_start': array([279], dtype=int32)}
 4 {'text': array(['a Marian place of prayer and reflection'], dtype=object), 'answer_start': array([381], dtype=int32)}
 5 {'text': array(['a golden statue of the Virgin Mary'], dtype=object), 'answer_start': array([92], dtype=int32)}
 6 {'text': array(['September 1876'], dtype=object), 'answer_start': array([248], dtype=int32)}
 7 {'text': array(['twice']), 'answer_start': array([441], dtype=int32)}
 8 {'text': array(['The Observer']), 'answer_start': array([598], dtype=int32)}
 9 {'text': array(['three']), 'answer_start': array([126], dtype=int32)}
10 {'text': array(['1987']), 'answer_start': array([908], dtype=int32)}
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb/data/squad$
```

# Project Details & Guidance (SQuAD Team)

The screenshot shows a red header with the arXiv logo and a navigation bar for 'Computer Science > Computation and Language'. Below the header, the title 'No Language Left Behind: Scaling Human-Centered Machine Translation' is displayed, along with the names of the NLLB Team members. There are buttons for 'CODE', 'Notebook', and social sharing. A text box contains a summary of the research, mentioning the goal of eradicating language barriers and the development of a conditional compute model based on Sparsely Gated Mixture of Experts. At the bottom, there's a call-to-action button 'Ask the author(s) a question! :)' and a 'powered by CatalyzeX' logo.

[Submitted on 11 Jul 2022 (v1), last revised 25 Aug 2022 (this version, v3)]

## No Language Left Behind: Scaling Human-Centered Machine Translation

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangtip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, Jeff Wang (NLLB Team)

Ask the author(s) a question! :)

Ask

powered by CatalyzeX

Driven by the goal of eradicating language barriers on a global scale, machine translation has solidified itself as a key focus of artificial intelligence research today. However, such efforts have coalesced around a small subset of languages, leaving behind the vast majority of mostly low-resource languages. What does it take to break the 200 language barrier while ensuring safe, high quality results, all while keeping ethical considerations in mind? In No Language Left Behind, we took on this challenge by first contextualizing the need for low-resource language translation support through exploratory interviews with native speakers. Then, we created datasets and models aimed at narrowing the performance gap between low and high-resource languages. More specifically, we developed a conditional compute model based on Sparsely Gated Mixture of Experts that is trained on data obtained with novel and effective data mining techniques tailored for low-resource languages. We propose multiple architectural and training improvements to counteract overfitting while training on thousands of tasks. Critically, we evaluated the performance of over 40,000 different translation directions using a human-translated benchmark, Flores-200, and combined human evaluation with a novel toxicity benchmark covering all languages in Flores-200 to assess translation safety. Our model achieves an improvement of 44% BLEU relative to the previous state-of-the-art, laying important groundwork towards realizing a universal translation system. Finally, we open source all contributions described in this work, accessible at [this https URL](https://this https URL).

- a conditional compute model based on Sparsely Gated Mixture of Experts

# Project Details & Guidance (SQuAD Team)

- NLLB အကြောင်းစိတ်ဝင်စားရင် စာတမ်းကို ဖတ်ပါ။ စာမျက်နှာ ၁၉၂ မျက်နှာ ရုံးပါတယ်။
- GitHub Link: <https://github.com/facebookresearch/fairseq/tree/nllb>
- ဆရာ ကိုယ်တိုင် NLLB ကို သုံးပြီး translation experiment တချို့ စမ်းလုပ်ထားတာကို log မှတ်ထားတပါတယ်။ အဲဒီ log ဖိုင်ကို လွှဲလာရင်လည်း အဆင့်ဆင့် လက်တွေ့လုပ်ဆောင်ရွက် အပိုင်းကို အကြမ်း နားလည်သွားပါလိမ့်မယ်။
- <https://github.com/ye-kyaw-thu/error-overflow/blob/master/nllb-exp.md>

# Project Details & Guidance (SQuAD Team)

- NLLB translation လုပ်ပေးတဲ့ shell script က ဆာဟပါမှာ ရှိပြီးသား
- squad\_train2my.sh, squad\_valid2my.sh နှစ်ဖိုင်ကို ပြင်ဆင်ခဲ့

```
#!/bin/bash

# Base directory for input files
INPUT_DIR="/home/ye/ye/exp/gpt-mt/nllb/data/squad/train/"

# Directory for output files
OUTPUT_DIR="/home/ye/ye/exp/gpt-mt/nllb/squad-my/"

# Create the output directory if it doesn't exist
mkdir -p "$OUTPUT_DIR"

# Iterate over each .src file in the input directory
for FILE in "$INPUT_DIR"/*.txt; do
    # Extract the base filename without the extension
    BASENAME=$(basename "$FILE" .txt)

    # Define the output file name
    OUTPUT_FILE="$OUTPUT_DIR/$BASENAME.my"

    # Print the command being executed (for debugging)
    echo "Running nllb-translate.sh for $FILE"

    # Run the translation command
    time ./nllb-translate.sh --input "$FILE" --source eng_Latn --target mya_Myrm --output "$OUTPUT_FILE"
done
```

# Project Details & Guidance (SQuAD Team)

- translation process က အနည်းဆုံး တပတ်လောက် ကြာနိုင်တာမိုလို server ကို ချိတ်ထားတဲ့ terminal ရဲ့ connection က ပြတ်သွားရင် program က ရပ်သွားနိုင်တာမိုလို screen command ကို သုံးကြရအောင်။
- screen command ကို သုံးတတ်မှ ဖြစ်မယ်။
- \$screen -S squad (screen တစ်ခု ဖန်တီးချင်ရင်)
- \$screen -ls (လက်ရှိ run ထားတဲ့ sceen တွေကို ရှာကြည့်ချင်ရင်)
- -D က run နေတဲ့ screen ကို detached လုပ်ဖို့ -R က reattached လုပ်ဖို့

```
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/nllb$ screen -ls
There are screens on:
    1955959.squad_valid      (06/12/2025 03:04:55 PM)          (Detached)
    1954893.squad            (06/12/2025 02:41:48 PM)          (Detached)
2 Sockets in /run/screen/S-ye.
```

# Project Details & Guidance (SQuAD Team)

- ဒီ ဆလိုက်ကို ပြင်နေတဲ့ အချိန်ထိ ဆာဗာပါမှာ ဘာသာပြန်ပြီးတဲ့  
စာကြောင်းအရေအတွက် က အောက်ပါအတိုင်း

```
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/n11b/squad-my$ ls  
train_context.my  valid_context.my  
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/n11b/squad-my$ date  
Fri Jun 13 03:31:21 PM +07 2025  
(base) ye@lst-gpu-server-197:~/ye/exp/gpt-mt/n11b/squad-my$ wc {train,valid}_context.my  
 5347    882407 12800456 train_context.my  
 5457    875469 12728282 valid_context.my  
10804   1757876 25528738 total
```

# Project Details & Guidance (SQuAD Team)

- Machine translation ကို တချို့က အဆင်ပြသူးပြလို့ ထင်နေကြပေမဲ့ လက်တွေ့မှာ ရှင်းရတဲ့ ပြဿနာတွေ အများကြီးပါ။ ဘာသာပြန် အမှားတွေ အများကြီးပါ
- တခါတလေ NLLB က ဘာသာပြန်မပေးနိုင်တဲ့ စာကြောင်းတွေလည်း ရှိတမိုလို trace လိုက်လို့ ရအောင် ဆရာ bash shell script မှာ parallel output (input ရော့ translated output ရော့) အနေနဲ့ဖိုင်မှာ သုမ္ပါဒေါ် ပြင်ထားပါတယ်။
- လက်ရှိ အချိန်ထိ ဘာသာပြန်ပြီးသလောက် စာကြောင်းတွေကို ဝင်လေ့လာကြည့်ကြရအောင်

# Project Details & Guidance (SQuAD Team)

- 19 The College of Engineering was established in 1920, however, early courses in civil and mechanical engineering were a part of the College of Science since the 1870s. Today the college, housed in the Fitzpatrick, Cushing, and Stinson-Remick Halls of Engineering, includes five departments of study – aerospace and mechanical engineering, chemical and biomolecular engineering, civil engineering and geological sciences, computer science and engineering, and electrical engineering – with eight B.S. degrees offered. Additionally, the college offers five-year dual degree programs with the Colleges of Arts and Letters and of Business awarding additional B.A. and Master of Business Administration (MBA) degrees, respectively.
- စက်မှုတက္ကသိလ်သည် ၁၉၂၀ ခုနှစ်တွင် တည်ထောင်ခဲ့သော်လည်း ၁၈၇၀ ခုနှစ်မှစ၍ အရပ်ဘက်နှင့် စက်မှုအင်ဂျင်နှီယာပညာရပ်များတွင် အကောင်းသင်တန်များသည် သိပ္ပါတက္ကန်၏ အစိတ်အပိုင်းတစ်ခုပြစ်သည်။ ယနေ့တွင် Fitzpatrick, Cushing နှင့် Stinson-Remick Halls of Engineering တွင် တည်ရှိသည့် ကောလိပ်တွင် ပညာရေးနာနဲ့ခုပါဝင်သည်။ လေကြောင်းနှင့် စက်ပစ္စည်းအင်ဂျင်နှီယာ၊ ဓာတုနှင့် ဒီဇိုင်းလီကျိုးအင်ဂျင်နှီယာ၊ အရပ်ဘက်နှင့် ဘူမိမေဒသပို့ ကုန်ပျိုးတာသိပိုင်း အင်ဂျင်နှီယာရှင်း လျှပ်စစ်အင်ဂျင်နှီယာ စ ခုဖြင့် BS ဘူမှုးပေးသည်။ ထိုပြင် ကောလိယပ်သည်အခြားအထူးတန်များဖြစ်သော B.A နှင့် Master of Business Administration (MBA) ဘူမှုးကိုပေးသည်။
- 20 The College of Engineering was established in 1920, however, early courses in civil and

- ဘာသာပြန်ပြီး ရလ္လာတဲ့ ဖိုင်တွေထဲက train context.my ကိုဖွင့်ပြီး လိုင်းနံပါတ် ၁၉ ရဲ့ input/output ကို လေ့လာလဲ ရအောင် တဲ့ပြုပေးထားတာပါ။
- သိပ္ပါတက္ကသိလ် ⇒ သိပ္ပါတက္ကစ် ဆုတဲ့ အများမျိုး၊ ကောလိယပ် ဆုတဲ့ အများမျိုးနဲ့ တချို့ကို အကောလိပ်စာကြောင်း အတူငါး ယူချလာတာမျိုးတွေ တွေ့ရပါလိမ့်မယ်

# Project Details & Guidance (SQuAD Team)

4314

The Amazon rainforest (Portuguese: Floresta Amazônica or Amazônia; Spanish: Selva Amazónica, Amazonía or usually Amazonia; French: Forêt amazonienne; Dutch: Amazoneregenwoud), also known in English as Amazonia or the Amazon Jungle, is a moist broadleaf forest that covers most of the Amazon basin of South America. This basin encompasses 7,000,000 square kilometres (2,700,000 sq mi), of which 5,500,000 square kilometres (2,100,000 sq mi) are covered by the rainforest. This region includes territory belonging to nine nations. The majority of the forest is contained within Brazil, with 60% of the rainforest, followed by Peru with 13%, Colombia with 10%, and with minor amounts in Venezuela, Ecuador, Bolivia, Guyana, Suriname and French Guiana. States or departments in four nations contain "Amazonas" in their names. The Amazon represents over half of the planet's remaining rainforests, and comprises the largest and most biodiverse tract of tropical rainforest in the world, with an estimated 390 billion individual trees divided into 16,000 species.

အမာဇာန်သစ်တော် (အမေရิกုံနှင့်တော်) သည် တောင်အမေရิกကနိုင်ရှိ အမာဇာန်ပင်လယ်အောင်ဒေသ၏ အများစုကို ဖုံးအပ်သည့် စိတ်ဝင်းသော သစ်တော်ဖြစ်သည်။ အမာဇာန် ပင်လယ်အိုသည် စတုရန်းကိုလိုမိတာ ၂၀၀၀၀၀၀ (၂၇၀၀၀၀၀ စတုရာမိုင်) ကို ဖုံးလွှမ်းထားသည်။ ယင်းအနက် ၅၅၀၀၀၀၀ (၁၀၀၀၀၀) စတုရာမိုင်သည် အမာဇာန်သစ်တော်ပြင် ဖုံးကွဲယ်ထားသည်။ ဤဒေသတွင် နိုင်ငံ ၉ နိုင်ငံပိုင်လျှပ်းသော နယ်မြေများပါဝင်သည်။ သစ်တော်အများစုသည် ဘရာမီးနိုင်အတွင်းရှိရာ၊ အမာရှတ်သစ်ပင် ၆၀ ရာခိုင်နှင့်၊ ပီရူးနိုင်မှ ၁၃၅%၊ ပင်နှီးလားနိုင်မှ ၁၀%၊ အီဂ္ဂါဒီ၊ ဘီလီးဒီးယား၊ ရှိုင်ယာ၊ ဆူရီန်နှင့် ပြင်သစ်ရှိုင်နားနိုင်များတွင် အနည်းငယ်ပါဝင်သည်။ ကောလ္ဗားနိုင်ရှိ နိုင်ငံ၊ ဤခုံမဟုတ် "အမာဇာန်" နယ်ပယ်သည်၍ ငါးတို့အမှည့်များတွင် ပါဝင်သည်။ အမာဇာ

4315

The Amazon rainforest (Portuguese: Floresta Amazônica or Amazônia; Spanish: Selva

- ဘာသာပြန်ပြီး ရလာတဲ့ ဖို့ငွေတွေထဲက valid\_context.my ထဲက ဥပမာပါ။  
အရမ်းရှည်တဲ့ စာကြောင်းတွေခုံ့ရင် NLLB က ဖြတ်ချေပစ်လိုက်တာကို မြင်ကြရပါလဲမှုမယ့်

# Project Details & Guidance (SQuAD Team)

- SQuAD Team က NLLB အပြင် တခြား NMT (e.g. Google-Translate) ဒါမှုမဟုတ် LLM (e.g. DeepSeek, ChatGPT) ကို သုံးပြီး SQuAD dataset တစ်ခုလုံးကို ဘာသာပြန်တာကို လက်တွေ့လှပါကြည့်စေချင်တယ်
- ဘာသာပြန်ပြီးထွက်လာတဲ့ အမှားတွေကို manual လှက်စစ်ပြီး ပြင်နိုင်သလောက် ပြင်တာမျိုး လုပ်ဖို့လုံးအပ်လိမ့်မယ်။ တကယ်တမ်း အသုံးဝင်တဲ့ Myanmar SQuAD dataset အနေနဲ့ အများသုံးလို့ ရအောင် ရှုချင်တယ် ဆိုရင်
- Modeling (ဥပမာ စာပိုဒ်ကို ဖတ်ခိုင်းပြီး မေးတာကို ဘယ်လောက်ဖြေနိုင်သလဲ experiment) ပါ လုပ်ကြည့်ကြဖို့အတွက် ပြင်ဆင်ဖို့ လုံးနေတာ နောက်တစ်ခုက ဘာသာပြန်ထားတဲ့ မြန်မာစာ ဒေတာကနေ မှန်တဲ့ အဖြေတွေက ဘာတွေဆုံးတာကို ကြိုတင် သတ်မှတ်ရန်

# Project Details & Guidance (SQuAD Team)

- answer field ይችል በformat ከ ማረኞቸው፡ ቁጥር ማረኞቸውን በመሆኑ፣ በምንጫል በመሆኑ

```
87581 {'text': array(['All Nepal Football Association'], dtype=object), 'answer_start': array([160], dtype=int32)}
87582 {'text': array(['25,000'], dtype=object), 'answer_start': array([498], dtype=int32)}
87583 {'text': array(['Tripureshwor'], dtype=object), 'answer_start': array([430], dtype=int32)}
87584 {'text': array(['Chinese'], dtype=object), 'answer_start': array([628], dtype=int32)}
87585 {'text': array(['17,182'], dtype=object), 'answer_start': array([54], dtype=int32)}
87586 {'text': array(['hilly terrain'], dtype=object), 'answer_start': array([289], dtype=int32)}
87587 {'text': array(['BP'], dtype=object), 'answer_start': array([500], dtype=int32)}
87588 {'text': array(['west'], dtype=object), 'answer_start': array([457], dtype=int32)}
87589 {'text': array(['Araniko'], dtype=object), 'answer_start': array([466], dtype=int32)}
87590 {'text': array(['Tribhuvan International Airport'], dtype=object), 'answer_start': array([71], dtype=int32)}
87591 {'text': array(['6'], dtype=object), 'answer_start': array([134], dtype=int32)}
87592 {'text': array(['22'], dtype=object), 'answer_start': array([297], dtype=int32)}
87593 {'text': array(['Amsterdam'], dtype=object), 'answer_start': array([698], dtype=int32)}
87594 {'text': array(['Turkish Airlines'], dtype=object), 'answer_start': array([734], dtype=int32)}
87595 {'text': array(['Oregon'], dtype=object), 'answer_start': array([229], dtype=int32)}
87596 {'text': array(['Rangoon'], dtype=object), 'answer_start': array([414], dtype=int32)}
87597 {'text': array(['Minsk'], dtype=object), 'answer_start': array([476], dtype=int32)}
87598 {'text': array(['1975'], dtype=object), 'answer_start': array([199], dtype=int32)}
87599 {'text': array(['Kathmandu Metropolitan City'], dtype=object), 'answer_start': array([0], dtype=int32)}
```

# Project Details & Guidance (SQuAD Team)

- SQuAD ပရောဂျက်ကို အခြေခံပြီး လက်တွေ့လုပ်ပြထားတဲ့ preprocessing, NLLB translation အပိုင်းက တွေ့ကြား  
ပရောဂျက်ခေါင်းစဉ်ကို လုပ်ဖို့ စဉ်းစားထားတဲ့ သူတွေအတွက်လည်း  
အသုံးဝင်ပါလိမ့်မယ်။
- Linux command, Python programming, shell script writing တွေကတော့  
ကွန်ပျူးတာကျောင်းသားတိုင်းက မသိမဖြစ်လို့ ပြောလို့ ရပါတယ်။  
အထူးသဖြင့်တော့ machine learning, NLP/AI အလုပ်တွေကို လုပ်ချင်တဲ့  
သူတွေအတွက်က
- Linux command အသုံးပြုပုံတွေနဲ့ တွေ့ကြား programming နဲ့  
ပတ်သက်တောက အချိန်သပ်သပ်ကြား ယူပြီး သင်ပေးဖို့ အစီအစဉ်  
မရှိပေမဲ့ ကြိုရင် ကြိုသလို အခုလုံ လုပ်ပြသွားပါမယ်

# Project Details & Guidance (ASR/TTS Team)

- PyQt6 module ကို installationလုပ်ဖို့လိုအပ်လိမ့်မယ်

```
(base) C:\Users\801680\.spyder-py3\intern3>python ./recorder.py --help
Traceback (most recent call last):
  File "C:\Users\801680\.spyder-py3\intern3\recorder.py", line 7, in <module>
    from PyQt6.QtWidgets import (
ModuleNotFoundError: No module named 'PyQt6'

(base) C:\Users\801680\.spyder-py3\intern3>pip install PyQt6
WARNING: Ignoring invalid distribution -cipy (c:\users\801680\anaconda3\lib\site-packages)
Collecting PyQt6
  Downloading pyqt6-6.9.1-cp39-abi3-win_amd64.whl (25.7 MB)
  ━━━━━━━━━━━━━━━━ 25.7/25.7 MB 880.7 kB/s eta 0:00:00
Collecting PyQt6-sip<14,>=13.8 (from PyQt6)
  Downloading pyqt6_sip-13.10.2-cp39-cp39-win_amd64.whl (53 kB)
  ━━━━━━━━━━━━━━ 53.6/53.6 kB 398.1 kB/s eta 0:00:00
Collecting PyQt6-Qt6<6.10.0,>=6.9.0 (from PyQt6)
  Downloading pyqt6_qt6-6.9.1-py3-none-win_amd64.whl (73.8 MB)
  ━━━━━━━━━━━━━━ 73.8/73.8 MB 812.2 kB/s eta 0:00:00
WARNING: Ignoring invalid distribution -cipy (c:\users\801680\anaconda3\lib\site-packages)
Installing collected packages: PyQt6-Qt6, PyQt6-sip, PyQt6
Successfully installed PyQt6-6.9.1 PyQt6-Qt6-6.9.1 PyQt6-sip-13.10.2

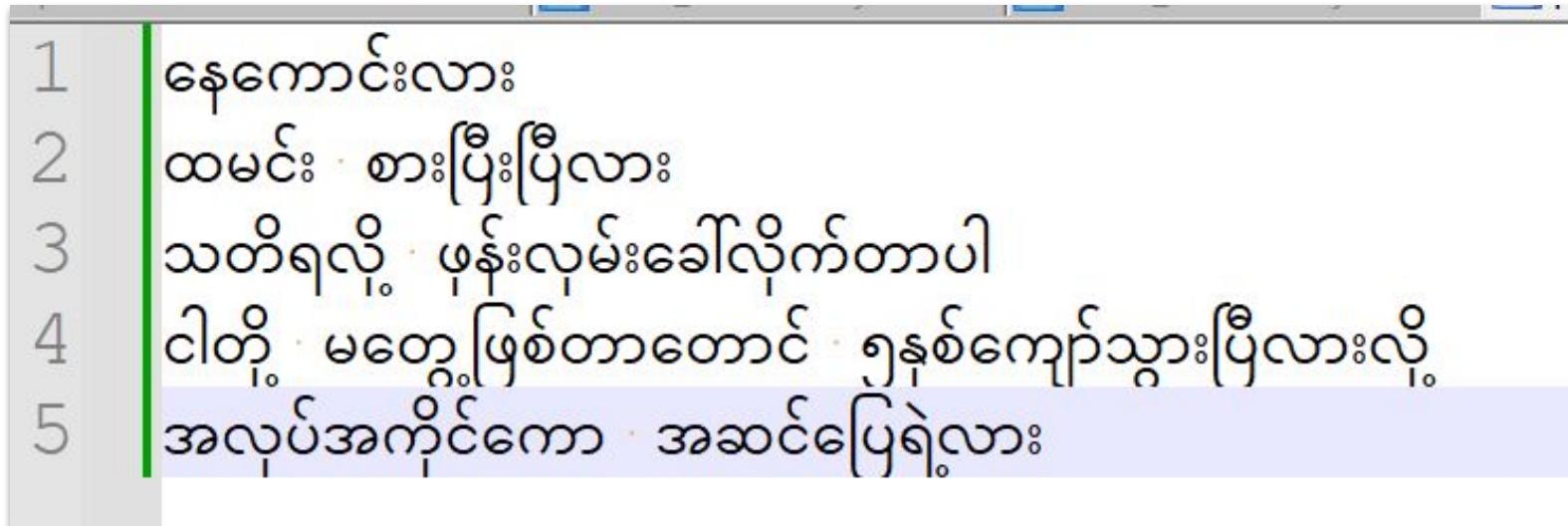
(base) C:\Users\801680\.spyder-py3\intern3>
```

# Project Details & Guidance (ASR/TTS Team)

- Domain ကို သတ်မှတ်ကြရအောင်
- Recording ကို ဘယ်လုပ်လုပ်ကြမလဲ ဆိုတာကို တိုင်ပင်ကြရလိမ့်မယ်
- ASR အတွက်ဆိုရင် ကုသု member တွေအပြင် တခြားသူတွေကိုပါ အကူအညတောင်းနှင်ရင် ပိုကောင်းတယ်
- Crowdsourcing လုပ်ကြမလား
- Mic ဘာညာမဝယ်နှင်ကြရင်တောင်၊ notebook မှာ built in ပါတဲ့ mic နဲ့ သွင်းကြတာပါ
- PC ကုပဲ သုံးပြီး Recording လုပ်ကြဖို့အတွက် recorder program ကိုတော့ ဆရာ ရေးပေးထားပါတယ်

# Project Details & Guidance (ASR/TTS Team)

- Example prompt file, prompt.txt



The screenshot shows a text editor window with a light gray background. On the left side, there is a vertical column of numbers from 1 to 5, each aligned with a specific line of text. A thick green vertical bar is positioned to the left of the first line (number 1). The text is in Burmese script. Lines 1 through 4 are in black font, while line 5 is in a lighter purple font.

1	နေကောင်းလား
2	ထမင်း စားပြီးပြီလား
3	သတိရလို့ ဖုန်းလုမ်းချေလိုက်တာပါ
4	ငါတို့ မတွေ့ဖြစ်တာတောင် ၅နှစ်ကျော်သွားပြီလားလို့
5	အလုပ်အကိုင်ကော် အဆင်ပြုရဲ့လား

# Project Details & Guidance (ASR/TTS Team)

- Usage examples
- `python recorder.py -p prompts.txt -d asr_recordings -m random -sr 16000`
- `python recorder.py -p chapter1.txt -d audiobook -m sequential -a`
- `python recorder.py -p phrases.txt -l 50 -m ordered -c 20`
- -d option မပေးရင် recording လုပ်တဲ့ အချိန် ရက်စွဲနဲ့ folder အသစ်တစ်ခုအောက်ပြီး သိမ်းပေးလိမ့်မယ် (e.g. `rec_0518_14Jun2025`)

# Project Details & Guidance (ASR/TTS Team)

- ဒီ recorder ကတော့ python နဲ့  
ရေးထားတာ
- Recording အတွက်က  
sounddevice library ကို သုံးခဲ့
- UI အတွက်က PyQt6 ကို သုံးခဲ့
- Mouse အပြင် long-term  
recording အတွက်  
အဆင်ပြေအောင် shortcut key ပါ  
ထည့်ပေးထား
- Recording  
လုပ်ခဲ့တဲ့အတွေအာကြံအရ<sup>၁</sup>  
အဆင်ပြေဆုံးဖြစ်အောင် minimal  
လုပ်ပေးထားပါတယ်

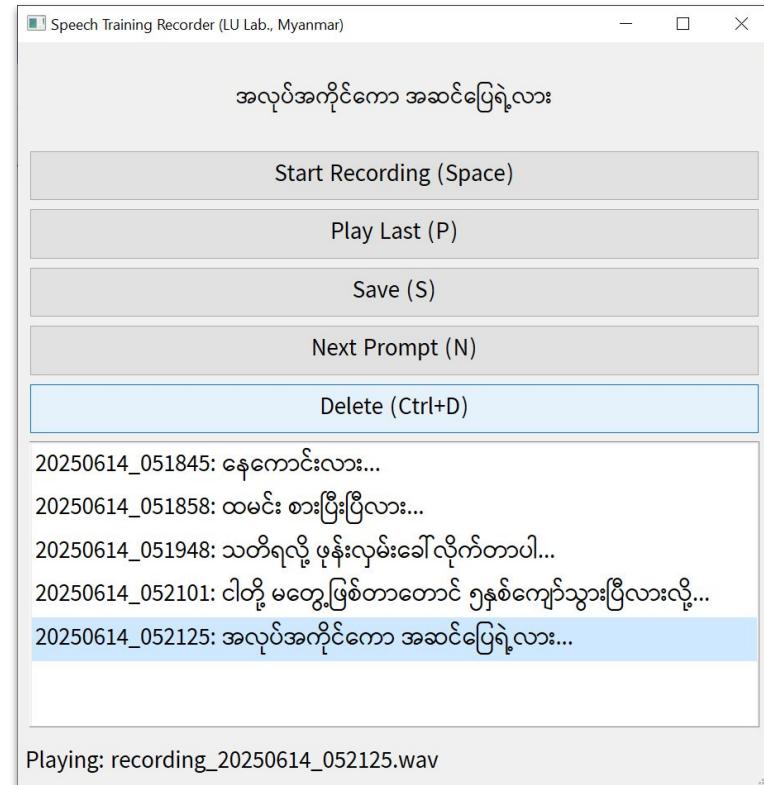


Fig. UI of Speech Training Recorder

# Project Details & Guidance (ASR/TTS Team)

- Random, Ordered, Sequential အလုပ်လုပ်ပုံက

[Prompt File]

"နေကောင်းလား"  
"ထမင်း စားပြီးပြီလား"  
"သတိရလို့ ဖုန်းလှမ်းခေါလိုက်တာပါ"

Mode	1st N	2nd N	3rd N	4th N	...
random	B	A	C	(stops)	
ordered	A	B	C	(stops)	
sequential	A	B	C	A	(loops)

# Project Details & Guidance (ASR/TTS Team)

- Example output folder

```
(base) C:\Users\801680\.spyder-py3\intern3>dir rec_0518_14Jun2025
Volume in drive C has no label.
Volume Serial Number is 9C54-A208

Directory of C:\Users\801680\.spyder-py3\intern3\rec_0518_14Jun2025

06/14/2025  05:21 AM    <DIR>          .
06/14/2025  05:21 AM    <DIR>          ..
06/14/2025  05:21 AM                738 recordings.tsv
06/14/2025  05:18 AM                44,140 recording_20250614_051845.wav
06/14/2025  05:18 AM                78,252 recording_20250614_051858.wav
06/14/2025  05:19 AM                88,236 recording_20250614_051948.wav
06/14/2025  05:21 AM                119,020 recording_20250614_052101.wav
06/14/2025  05:21 AM                75,756 recording_20250614_052125.wav
                           6 File(s)   406,142 bytes
                           2 Dir(s)  29,852,176,384 bytes free

(base) C:\Users\801680\.spyder-py3\intern3>
```

# Project Details & Guidance (ASR/TTS Team)

- နောက်ထပ် အရေးကြီးတဲ့ output ဖိုင်က recordings.tsv ဖိုင်ပါ
- Column or field ၅ခုခုခုပြီး သိမ်းပေးလိမ့်မယ်
- Delimiter ၏ TAB
- filename, prompt, timestamp, sample\_rate, bit\_depth
- Example recordings.tsv file

	filename	prompt	timestamp	sample_rate	bit_depth
1	recording_20250614_051845.wav	နောက်ငါးလား	20250614_051845	16000	16
2	recording_20250614_051858.wav	ထာမင်း စားပြီးပြီလား	20250614_051858	16000	16
3	recording_20250614_051948.wav	သတိရလို့ ဖုန်းလှမ်းခေါ်လိုက်တာပါ	20250614_051948	16000	16
4	recording_20250614_052101.wav	ဝါတို့ မတွေ့ဖြစ်တောင် ရှန်ကျော်သွားပြီလားလို့	20250614_052101	16000	16
5	recording_20250614_052125.wav	အလုပ်အကိုင်ကော် အဆင်ပြုရဲလား	20250614_052125	16000	16
6					
7					

# Project Details & Guidance (ASR/TTS Team)

- Linux မှာဆိုရင် sox, soxi ဆိုတဲ့ အသုံးဝင်တဲ့ command ရှိတယ်

```
(base) ye@lst-gpu-server-197:~/intern3/rec_0518_14Jun2025$ soxi ./recording_20250614_052125.wav
```

```
Input File      : './recording_20250614_052125.wav'  
Channels       : 1  
Sample Rate    : 16000  
Precision      : 16-bit  
Duration       : 00:00:02.37 = 37856 samples ~ 177.45 CDDA sectors  
File Size      : 75.8k  
Bit Rate       : 256k  
Sample Encoding: 16-bit Signed Integer PCM
```

```
(base) ye@lst-gpu-server-197:~/intern3/rec_0518_14Jun2025$ -
```

# Project Details & Guidance (ASR/TTS Team)

- အသုံးပြုပုံ အသေးစိတ်ကို command line help ခေါ်ကြည့်ပါ

```
(base) C:\Users\801680\.spyder-py3\intern3>python ./recorder.py --help
usage: recorder.py [-h] [-p PROMPTS_FILENAME] [-d SAVE_DIR] [-m {random, ordered, sequential}]
                   [-c PROMPTS_COUNT] [-l PROMPT_LEN_SOFT_MAX] [-a] [-sr {8000, 16000, 44100, 48000}]
                   [-b {16, 32}]

Speech Training Recorder (LU Lab., Myanmar) – Record prompted speech

optional arguments:
  -h, --help            show this help message and exit
  -p PROMPTS_FILENAME, --prompts_filename PROMPTS_FILENAME
                        text file containing prompts (one per line)
  -d SAVE_DIR, --save_dir SAVE_DIR
                        custom output directory (default: auto-generated)
  -m {random, ordered, sequential}, --prompt_selection {random, ordered, sequential}
                        prompt selection mode (default: random)
  -c PROMPTS_COUNT, --prompts_count PROMPTS_COUNT
                        max prompts to use (default: 100)
  -l PROMPT_LEN_SOFT_MAX, --prompt_len_soft_max PROMPT_LEN_SOFT_MAX
                        maximum prompt length in characters (0=no limit)
  -a, --auto_next        auto-advance to next prompt after save
  -sr {8000, 16000, 44100, 48000}, --sample_rate {8000, 16000, 44100, 48000}
                        sample rate in Hz (default: 16000)
  -b {16, 32}, --bit_depth {16, 32}
                        bit depth (16 or 32, default: 16)

Example usages:
recorder.py -p prompts.txt -m random
recorder.py -p script.txt -m ordered -d custom_folder
recorder.py -p phrases.txt -m sequential -a
```

# Project Details & Guidance (ASR/TTS Team)

- ASR/TTS team အတွက်က ဘာရည်ရှယ်ချက်အတွက်လုပ်ကြမှုလဲ၊ ဘယ်လို အသုံးဝင်တဲ့ ပရောဂျက်ကို လုပ်ကြမှုလဲ ဆုံးတာကို စဉ်းစားကြရလိမ့်မယ်
- ဒုမ္မန်းကို သတ်မှတ်ပြီးတာနဲ့ text corpus ကို ပြင်ပြီး (i.e. prompt file) အသုံစသွင်းနိုင်အောင် လုပ်ကြရအောင်
- အသံက ၂လ လောက်တော့ သွေးလို ရလိမ့်မယ်
- Recording လည်း လုပ်ရွင်းနဲ့ တော်းလေ့လာစရာရှိတာတွေကို လေ့လာသွားကြရအောင်
- ရှေ့မှုလည်း ပြောခဲ့သလိုပဲ တတ်နိုင်သရှိ speech processing R&D pipeline ထက အရေးကြီးတဲ့ အပိုင်းတချို့ကို သံဃားလအတွင်း အတွေ့အကြိုရဖို့က အားလုံး ကြိုးစားကြမှ ဖြစ်လိမ့်မယ်

# Project Details & Guidance (Readability Team)

- Readability Score အတွက် အနည်းဆုံး evaluation dataset ကို  
ပြင်ကြရလိမ့်မယ်
- Level တွေ ဘယ်လိုခဲ့မလဲ၊ အဲဒီ ခွဲထားတဲ့ level တွေအတွက်  
လိုက်လျော့သီထွေဖြစ်တဲ့ ဒေတာကို စဉ်းစားကြရလိမ့်မယ်
- ပထမဆုံး seminar မှာ အကြမ်းရှင်းပြခဲ့သလုပ်ပါ။  
ဘာသာစကားတွေအတွက် readability score ကို ဘယ်လို  
သတ်မှတ်ကြသလဲ ဆုံးတာကို literature review သေချာလုပ်ကြရလိမ့်မယ်။  
မှတ်သားရလိမ့်မယ်
- မြန်မာစာအတွက် UFL က သတ်မှတ်ထားတဲ့ proficiency test  
ဆုံးတာကိုရှာဖွေလို့ ရနိုင်မလား
- Myanmar Language Test ကလည်း တိုက်ရှိက် မသက်ဆိုင်ပေမဲ့ သူတို့  
ဘယ်လို သတ်မှတ်ထားသလဲ လေ့လာသင့်တယ်

# Project Details & Guidance (Readability Team)

- Homepage of MLT
- Link: <https://www.mlt-myanmar.com/>

 **MLT** Myanmar Language Test

HOME About MLT For Test Takers Contact FAQ



# Myanmar Language Test



Test schedule					About MLT
No.	Test Date	Available Level	Place	Application Period	
28th MLT	Mar 16th, 2025 (Sun)	MB, M1	Online	Jan 20th, 2025 (Mon) ~ <b>Mar 9th, 2025 (Sun)</b>	<a href="#">Test Description</a>
29th MLT	Jun 14th, 2025 (Sat) Jun 15th, 2025 (Sun)	M2, M3 MB, M1	Online Online	Mar 17th, 2025 (Mon) ~ <b>Jun 8th, 2025 (Sun)</b>	<a href="#">Advantage</a>
30th MLT	Sep 21st, 2025 (Sun)	MB, M1	Online	Jun 16th, 2025 (Mon) ~ <b>Sep 7th, 2025 (Sun)</b>	<a href="#">Level (MB-M4) description</a>
31st MLT	Dec 20th, 2025 (Sat) Dec 21st, 2025 (Sun)	M2, M3, M4 MB, M1	Online Online	Sep 22nd, 2025 (Mon) ~ <b>Dec 7th, 2025 (Sun)</b>	<a href="#">Sample questions</a>

**About MLT**

- [Test Description](#)
- [Advantage](#)
- [Level \(MB-M4\) description](#)
- [Sample questions](#)
- [Reference / Books](#)
- [Partners](#)
- [Background / Objective](#)
- [Supervisor / Members](#)

**For Test Takers**

- [How to register](#)
- [Test Schedule](#)

# Project Details & Guidance (Readability Team)

Burmese alphabet and very basic vocabulary and grammar in daily life.

**MB**

Spoken / Written	Spoken: 100 %
Vocabulary range	200 – 300
Learning period	72 hours (3 months)

[MB guideline](#)

[MB Sample Question](#)

[MB Reference / Books](#)

Basic vocabulary and grammar to survive and communicate in daily life.

**M1**

Spoken / Written	Spoken: 100 %
Vocabulary range	500 – 700
Learning period	144 hours (6 months)

[M1 guideline](#)

[M1 Sample Question](#)

[M1 Reference / Books](#)

Be able to communicate in daily life and read simple paragraph.

**M2**

Spoken / Written	Spoken: 100 %
Vocabulary range	1200 – 1400
Learning period	288 hours (12 months)

[M2 guideline](#)

[M2 Sample Question](#)

[M2 Reference / Books](#)

Be able to communicate with native speaker and read normal paragraph.

**M3**

Spoken / Written	Spoken: 70 %, Written: 30 %
Vocabulary range	2500 – 3500
Learning period	432 hours (18 months)

[M3 guideline](#)

[M3 Sample Question](#)

[M3 Reference / Books](#)

Be able to communicate in Business situation and read difficult paragraph.

**M4**

Spoken / Written	Spoken: 50 %, Written: 50 %
Vocabulary range	6000 ++
Learning period	576 hours (24 months ++)

[M4 guideline](#)

[M4 Sample Question](#)

[M4 Reference / Books](#)

# Project Details & Guidance (Readability Team)

- MB:  
[https://www.mlt-myanmar.com/wp-content/uploads/2022/01/Guideline-for-the-test-taker-MB\\_202201.pdf](https://www.mlt-myanmar.com/wp-content/uploads/2022/01/Guideline-for-the-test-taker-MB_202201.pdf)
- M1:  
[https://www.mlt-myanmar.com/wp-content/uploads/2022/01/Guideline-for-the-test-taker-M1\\_202201.pdf](https://www.mlt-myanmar.com/wp-content/uploads/2022/01/Guideline-for-the-test-taker-M1_202201.pdf)
- M2:  
[https://www.mlt-myanmar.com/wp-content/uploads/2023/02/Guideline-for-the-test-taker-M2\\_202301.pdf](https://www.mlt-myanmar.com/wp-content/uploads/2023/02/Guideline-for-the-test-taker-M2_202301.pdf)
- M3:  
[https://www.mlt-myanmar.com/wp-content/uploads/2022/01/Guideline-for-the-test-taker-M3\\_202201.pdf](https://www.mlt-myanmar.com/wp-content/uploads/2022/01/Guideline-for-the-test-taker-M3_202201.pdf)
- M4:  
[https://www.mlt-myanmar.com/wp-content/uploads/2022/01/Guideline-for-the-test-taker-M4\\_202201.pdf](https://www.mlt-myanmar.com/wp-content/uploads/2022/01/Guideline-for-the-test-taker-M4_202201.pdf)

# Project Details & Guidance (Readability Team)

- I read this The Stanford Foreign Language Oral Skills Evaluation Matrix (FLOSEM): A Rating Scale for Assessing Communicative
- သုက Oral
- proficiency measures ကို အခိုကထားတယ်
- ပြီးတော့ foreign language အနေနဲ့စဉ်းတားထားတယ်
- ငါတိုက text ကိုပဲ အခိုကထားရမယ်

DOCUMENT RESUME	
ED 445 538	FL 026 412
AUTHOR	Padilla, Amado M.; Sung, Hyekyung
TITLE	The Stanford Foreign Language Oral Skills Evaluation Matrix (FLOSEM): A Rating Scale for Assessing Communicative Proficiency.
PUB DATE	1999-05-00
NOTE	38p.
PUB TYPE	Reports - Research (143)
EDRS PRICE	MF01/PC02 Plus Postage.
DESCRIPTORS	Chinese; *Communicative Competence (Languages); *Evaluation Methods; High School Students; High Schools; Interviews; Japanese; Korean; *Language Proficiency; *Oral Language; Rating Scales; Second Language Instruction; Second Language Learning; Speech Skills; *Student Evaluation; Uncommonly Taught Languages
ABSTRACT	
The development of the Stanford Foreign Language Oral Skills Evaluation Matrix (FLOSEM), a rating scale for assessing communicative proficiency in foreign languages, is described. Information on the utility of the FLOSEM is presented based on the results of three studies. Oral proficiency measures were obtained by means of the FLOSEM from 573 high school students enrolled in beginning through advanced Japanese, Chinese, and Korean. Classroom foreign language teachers rated students' proficiency at the beginning and the end of the school year. Students also used the FLOSEM to rate their own proficiency in the target language. In addition to FLOSEM ratings, oral proficiency was also assessed for a subset of 132 students by means of the Classroom Oral Competency Interview (COCI), which is a brief 5-7	

# Project Details & Guidance (Readability Team)



UPPSALA  
UNIVERSITET

Department of Linguistics and Philology  
Språkteknologiprogrammet  
(Language Technology Programme)  
Master's thesis in Computational Linguistics

14 juni 2006

- for Swedish
- By Patrik Larsson (2006)
- Supervised by Beata Megyesi
- ଶୀଘର ପାଠ୍ୟକର୍ତ୍ତା

## Classification into Readability Levels

Implementation and Evaluation

Link: <http://www.diva-portal.org/smash/get/diva2:131028/FULLTEXT01.pdf>

# Project Details & Guidance (Readability Team)

- ଟାଖୁର୍କଣ୍ଠା ଓ ମୂଳ

1. **Lorges Formula** (1939) revised in 1948, classifies texts into grades 3-12.

$$Grade = .07sl + .1073w_d + .1301pp + 1.6126 \text{ where}$$

- $sl$  = Average sentence length.
- $w_d$  = Number of different difficult words per 100 words. Difficult words are words not on the *Dale 769-word list*.
- $pp$  = Number of prepositional phrases per 100 words.

Lorges formula is considered as the best among the earliest formulas. It is relatively straightforward to calculate, which made it popular.

# Project Details & Guidance (Readability Team)

- ତାଖୁର୍କଣ୍ଡା ଓ ମୁଦ୍ରା
2. **Dale-Challs Formula** (1948) revised in 1995, classifies texts into grades 3-12.

$$Grade = .0596sl + .1579w_d + 3.6365$$

- $sl$  = Average sentence length.
- $w_d$  = The percentage of words not occurring on the *Dale list of 3000*.

In 1995, the *Dale list of 3000* was updated and the formula was changed. The reason why this formula is less known and used than, for example Flesch's Reading Ease formula is that it is more difficult to calculate, since checking of the 3000 words on the list is a time consuming task.

# Project Details & Guidance (Readability Team)

- Dale-Chall 3000 words (დალ-ტოლის 3000 woordი)

a able aboard about above absent accept accident account ache aching acorn acre across act acts add address admire adventure afar → afraid after afternoon afterward afterwards again against age aged ago agree ah ahead aid aim air airfield airplane airport airship airy → alarm alike alive all alley alligator allow almost alone along aloud already also always am America American among amount an and → angel anger angry animal another answer ant any anybody anyhow anyone anything anyway anywhere apart apartment ape apiece appear apple April → apron are aren't arise arithmetic arm armful army arose around arrange arrive arrived arrow art artist as ash ashes aside ask → asleep at ate attack attend attention August aunt author auto automobile autumn avenue awake awaken away awful awfully awhile ax axebaa babe babies back background backward backwards bacon bad badge badly bag bake baker bakery baking ball balloon banana band bandage bang banjo bank banker bar barber bare barefoot barely bark barn barrel base baseball basement basket bat batch bath bathe bathing bathroom bathtub → battle battleship bay be beach bead beam bean bear beard beast beat beating beautiful beautify beauty became because become becoming bed bedbug bedroom bedspread bedtime bee beech beef beefsteak beehive been beer beet before beg began beggar begged begin beginning begun behave behind being → believe bell belong below belt bench bend beneath bent berries berry beside besides best bet better between bib bible bicycle bid big bigger bill billboard

# Project Details & Guidance (Readability Team)

- Dale-Chall 3000 words (ଡିଲ୍ ଅଳ୍ପଃପୁଣ୍ଡଃ)

understand · underwear → undress · unfair · unfinished · unfold · unfriendly · unhappy →  
unhurt · uniform · United · States · unkind · unknown → unless · unpleasant · until ·  
unwilling · up · upon → upper · upset · upside · upstairs · uptown · upward → us · use · used  
useful · valentine · valley · valuable → value · vase · vegetable → velvet · very · vessel · →  
victory · view · village → vine · violet · visit → visitor · voice · votewag · wagon · waist ·  
wait · wake · waken · walk · wall · walnut · want · war · warm · warn · was · wash · washer · washtub  
wasn't · waste · watch · watchman · water · watermelon · waterproof · wave · wax → way · wayside ·  
we · weak · weakness · weaken · wealth · weapon · wear · weary · weather · weave · web · we'd · wedding ·  
Wednesday · wee · weed · week · we'll · weep · weigh · welcome · well · went · were → we're · west ·  
western · wet · we've · whale · what · what's · wheat · wheel · when · whenever · where · which · while ·  
whip · whipped · whirl · whisky · whiskey · whisper · whistle · white · who · who'd · whole · who'll ·  
whom · who's · whose · why · wicked · wide · wife · wiggle · wild · wildcat · will · willing · willow ·  
win · wind · windy · windmill · window · wine · wing · wink · winner · winter · wipe · wire · wise ·  
wish · wit · witch · with · without · woke · wolf · woman · women · won · wonder · wonderful · won't ·  
wood · wooden · woodpecker · woods · wool · woolen · word · wore · work · worker · workman · world · →  
worm · worn · worry · worse · worst · worth · would · wouldn't · wound · wove · wrap · wrapped · wreck ·  
wren · wring · write · writing · written · wrong · wrote · wrungyard · yarn · year · yell → yellow ·  
yes · yesterday · yet → yolk · yonder · you · you'd · you'll · young · youngster · your · →  
yours · you're · yourself · yourselves → youth · you've

# Project Details & Guidance (Readability Team)

3. **Flesch Reading Ease formula** (1948) revised several times, returns a score where a higher score indicates a more difficult text.

$$\text{ReadingEase} = 206.835 - 1.015sl - .846wl$$

- $sl$  = Average number of word per sentence.
- $wl$  = Number of syllables per 100 words.

*Flesch Reading Ease formula* returns a number between 0-100, where a higher score indicates that the text is harder to read. The formula is very simple to calculate, since the text passage to analyze has to be only 100 words and only two features need to be investigated. Reading Ease has become U.S governmental standard and most states in the U.S, require insurance forms to score at a certain level (around 40-50) to be valid. The formula is also used in several word processors as a service to test the documents readability. Flesch Reading Ease formula is the most used and well-known formula and it has influenced formulas for other languages because of its high correlation score and the simple calculation.

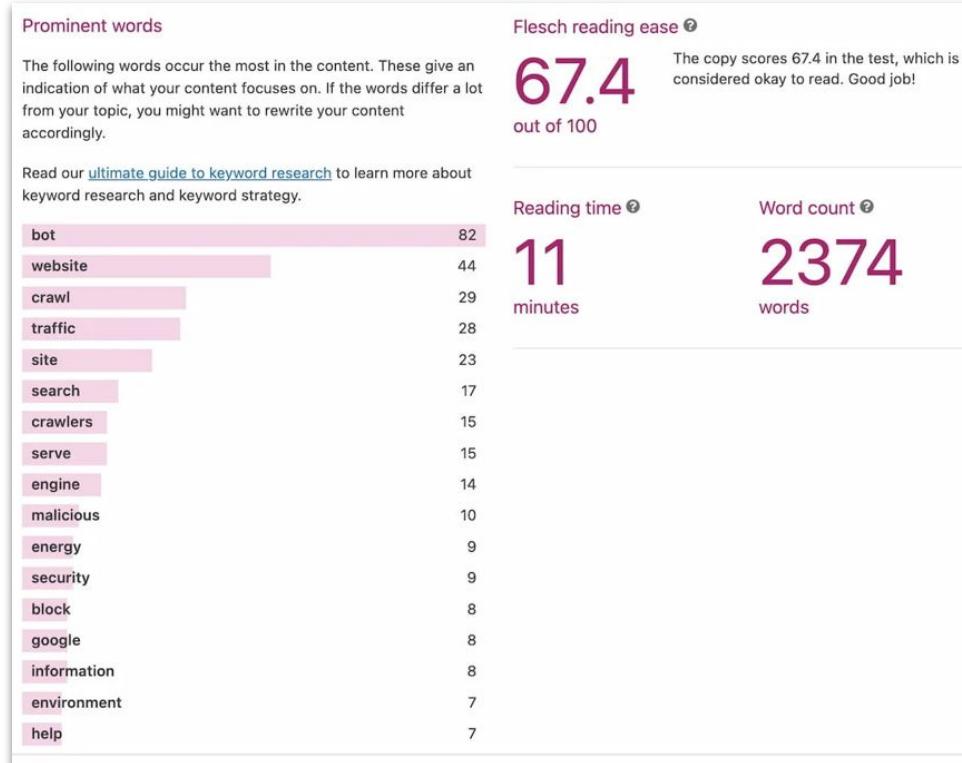


# Project Details & Guidance (Readability Team)

- the Flesch reading ease score
- Source: [https://en.wikipedia.org/wiki/Flesch%20%93Kincaid\\_readability\\_tests](https://en.wikipedia.org/wiki/Flesch%20%93Kincaid_readability_tests)

Score	School level (US)	Notes
100.00–90.00	5th grade	Very easy to read. Easily understood by an average 11-year-old student.
90.0–80.0	6th grade	Easy to read. Conversational English for consumers.
80.0–70.0	7th grade	Fairly easy to read.
70.0–60.0	8th & 9th grade	Plain English. Easily understood by 13- to 15-year-old students.
60.0–50.0	10th to 12th grade	Fairly difficult to read.
50.0–30.0	College	Difficult to read.
30.0–10.0	College graduate	Very difficult to read. Best understood by university graduates.
10.0–0.0	Professional	Extremely difficult to read. Best understood by university graduates.

# Project Details & Guidance (Readability Team)



- Reading time က ဘယ်လောက်ကြာသလဲ ဆိုတဲ့ အချက်ကို ခန့်မွှန်းတာကဗာလည်း စိတ်ဝင်စားဖို့ကောင်းတယ်

- Source: <https://yoast.com/yoast-seo-july-12-2022/>

# Project Details & Guidance (Readability Team)

- ତାମ୍ରିକା କାନ୍ତି

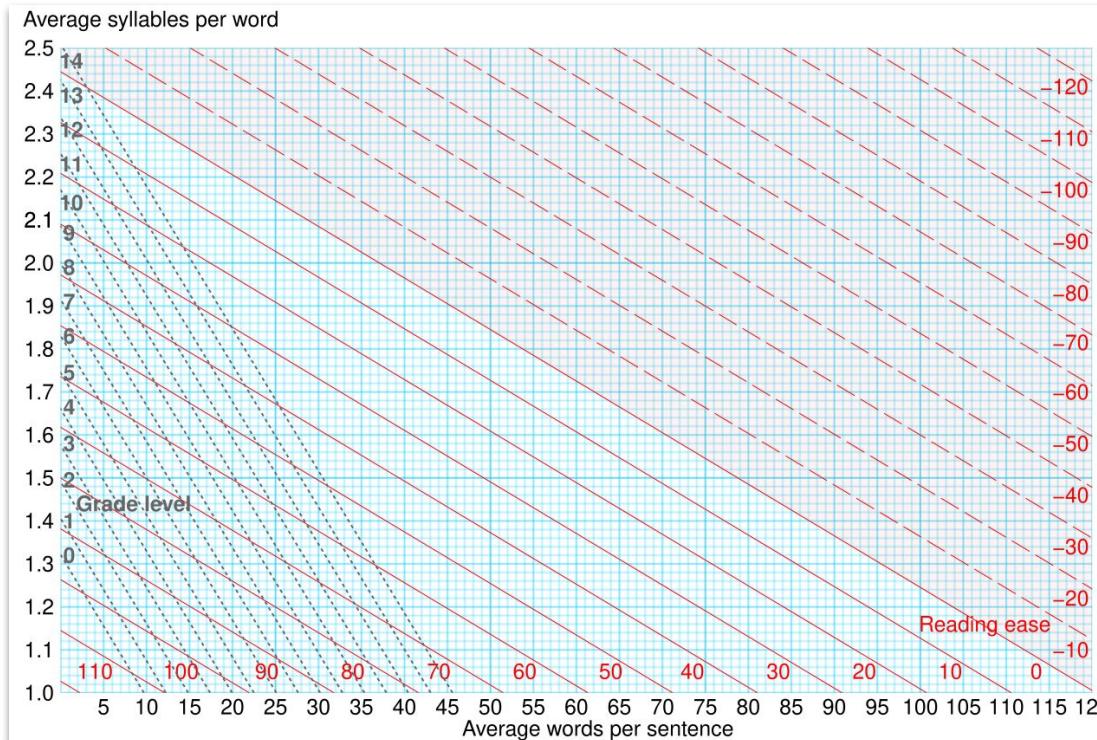
- 4. **Flesch-Kincaid Grade Level** (1975), classifies texts into American school grades.

$$Grade = .39sl + 11.8wl - 15.59$$

- $sl$  = Average number of word per sentence.
- $wl$  = Average number of syllables per word.

*Flesch-Kincaid Grade Level* is a modification of Fleschs *Reading Ease Formula*. It translates the former formula to an U.S grade level. A score of 10.2 indicates that the text is understandable for a 10th grade American student. Flesch-Kincaid Grade Level formula is the most used formula that classifies texts into a grade level.

# Project Details & Guidance (Readability Team)



- Graphs of Flesch-Kincaid reading ease (red) and grade level (gray) scores against average syllables per word and average words per sentence

# Project Details & Guidance (Readability Team)

- How to use the Flesch reading ease score to improve your writing?
  - အလွယ်ပြောရရင် စာကြောင်းတွေကို တံတိရေး
  - နားလည်ရ ခက်တဲ့ စာလုံးတွေကို တတ်နှင့်သရုံးလွှာ့သုံး
- Flesch Reading-Ease ရေား Flesch–Kincaid Grade Level ရောက score တွက်တာက တူတယ် မတူတာက weighting factors

# Project Details & Guidance (Readability Team)

- ଟାଖୁର୍କଣ୍ଡା ଓ ମୁଦ୍ରା

5. **LIX (Läsbartehetsindex)** (1968) revised a few times, developed for Swedish.

$$LIX = wl/s + 100 * w_d/wl$$

- $wl$  = Number of words in the text.
- $s$  = Number of sentences in the text.
- $w_d$  = Number of difficult words in the text, where difficult words is defined as words consisting of more than six letters.

The value from *LIX* has to be interpreted with a *LIX-interpreter*. There are several available, two examples of interpreters is represented in Figure 2.1 [Björnsson, 1971]. Depending on which *LIX-interpreter* is used, *LIX* can be applied to any level of text, just by adjusting the scale. *LIX* have successfully been applied to several other languages with the good results, by simply adjusting the scale of the interpreter [Cedergren, 1992].

# Project Details & Guidance (Readability Team)

- စာမျက်နှာ ၁၂ မှာ

LIX value	Description	LIX value	Description
20	Very easy	20-25	Children's books
30	Easy	31-35	Fiction
40	Average	40-45	Newspapers
50	Hard	50-55	Science reports
60	Very Hard	60-	Government texts, law texts ...

**Figure 2.1:** Two LIX-interpreters, the left is a simple version and the right an adapted interpreter.

- ရွှေသွေး၊ တေဇာ၊ ဝထ္ခ၊ သတင်းစာ၊ သိပုံနှင့်နည်းပညာ၊ ဥပဒေ၊ စတာတွေနဲ့ခဲ့နိုင်

# Project Details & Guidance (Readability Team)

For “Readability,” Amazon lists [Fog Index](#) and two indices developed by Rudolf Flesch under commission by the Navy. Keeping it brief, the Fog Index estimates the number of years of formal education required to understand the book. The [Flesch-Kincaid Index](#) similarly measures the U.S. grade level likely needed to understand the book, meaning Fog and Flesh-Kincaid numbers should be, in a perfect world, quite similar. The regular Flesch Index is based on a 100 point scale, with 100 being the easiest to read (for frame of reference, a college degree is considered necessary to read a book with a score of 0 to 30, while a 5th grader should be able to understand a book with a score of 90 to 100).

Got all that? Here’s what our sample books scored (Fog, Flesch-Kincaid, Flesch):

*Finnegans Wake*: 11.8, 9.3, 57.8

*Where I’m Calling From*: 5.7, 3.9, 84.2

*The Great Gatsby*: 14.4, 11.8, 48.8

*The Memory Keeper’s Daughter*: 7.5, 5.6, 76.7

*The Tipping Point*: 12.6, 10.1, 55.7

- လူတွေ  
သပ်ဂရုမပြုမိက  
တဲ့ Amazon ရဲ့  
Text Stats (i.e.  
Readability+)
- အခုန္ဓာက်ပုံင်း  
Amazon website  
မှာ မထွေထော့ဘူး။

# Project Details & Guidance (Readability Team)

## Text Stats

These statistics are computed from the text of this book. ([learn more](#))

Readability ( <a href="#">learn more</a> )		Compared with other books	
Fog Index:	8.4	18% are easier	82% are harder
Flesch Index:	69.7	18% are easier	82% are harder
Flesch-Kincaid Index:	6.1	16% are easier	84% are harder
Complexity ( <a href="#">learn more</a> )			
Complex Words:	11%	32% have fewer	68% have more
Syllables per Word:	1.5	30% have fewer	70% have more
Words per Sentence:	10.4	16% have fewer	84% have more
Number of			
Characters:	600,727	72% have fewer	28% have more
Words:	105,436	77% have fewer	23% have more
Sentences:	10,102	90% have fewer	10% have more
Fun stats			
Words per Dollar:	10,737		
Words per Ounce:	10,983		

- လျှောက်ညံပြီး ငါတို့  
မြန်မာစာအတွက်  
Readability Score ကို  
ဘယ်လိုတွက်ကျမလဲ  
ဆိုတာကို  
စဉ်းစားစေချင်တယ်
- Team member  
တွေအကြားမှာ တိုင်ပင်  
စေချင်တယ်

# Project Details & Guidance (Readability Team)

12

IEEE TRANSACTIONS ON PROFESSIONAL COMMUNICATION, VOL. PC-30, NO.1, MARCH 1987

## Readability Formulas: Useful or Useless?

GLEND A M. McCLURE

**Abstract – This interview with Dr. J. Peter Kincaid explores the value of readability formulas and computer editing systems to the engineer who writes on the job. Dr. Kincaid developed the Kincaid Readability Formula, the standard used in judging the reading levels of Department of Defense manuals. Currently a research psychologist for the Army Research Institute, Dr. Kincaid is working on automating the writing and delivery of technical information. Previously, he was a team leader with the Navy's Training Analysis and Evaluation Group [1].**

Among other things which we can reasonably measure: the number of commonly understood words, sentence complexity, the number of abstract ideas, and the use of personal pronouns. Beyond these factors, it takes the expertise of the writer and editor to judge organization of the text and whether or not the text conveys the proper information.

*How accurate is the Kincaid formula in determining the correct reading grade level of text?*

As a rule of thumb, it is accurate to within plus or minus one reading grade level.

# Project Details & Guidance (Readability Team)

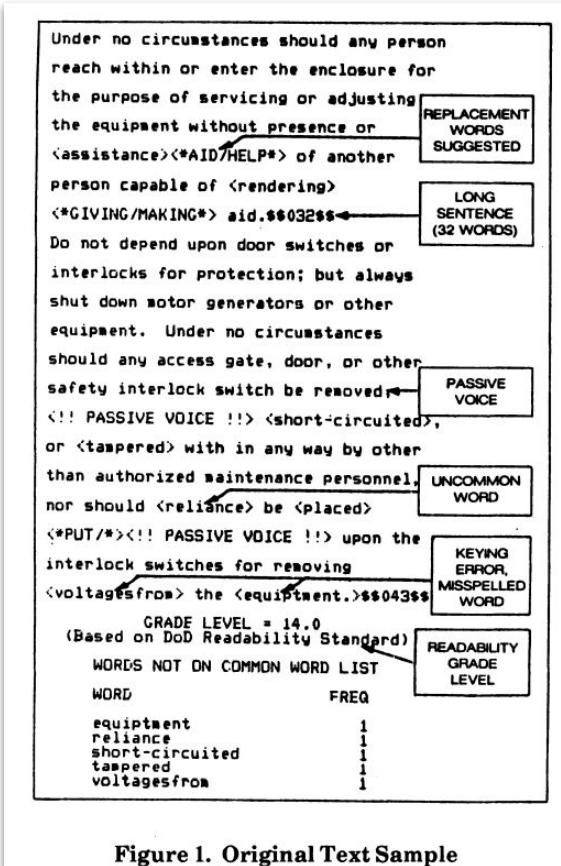


Figure 1. Original Text Sample

Do not reach within or enter the enclosure to service or adjust the equipment by yourself. Make sure another person able to help is with you. Do not depend upon door switches or interlocks for protection; always shut down motor generators or other equipment. Do not remove or short circuit any access gate, door, or other safety interlock switch. Only authorized maintenance personnel can do this. Do not depend on the interlock switches for removing voltages from the equipment.

GRADE LEVEL = 7.0  
(Based on DoD Readability Standard)

Figure 2. Rewritten Text Sample  
licensed use limited to: King Mongkut's Institute of Technology Ladkrabang provided by

# Project Details & Guidance (Readability Team)

- စဉ်းစားမိတာက ပထမဆုံး အဆင့်အနေနဲ့ readability ကို user-study or ဆာမေးလုပ်မယ်
- ဥပမာ ရွှေသွေး၊ တော့ ဝတ္ထု၊ သတင်းစာ၊ သံပံ့နှင့်နည်းပညာ၊ ဥပဒေ စာအုပ်တွေထဲကနေ စာကြောင်း၊ စာပိုဒ်တွေကို ရွှေးထုတ်ထားပြီး ကွဲကိုလပ်ဖြည့်ခိုင်းတာမျိုး လုပ်ခိုင်းပြီး ရလဒ်ကို ကြည့်မယ်
- ဒုတိယ တဆင့်အနေနဲ့က data driven approach ဖြစ်တဲ့ readability corpus ကို ဆောက်ထားပြီးမှ model ကနေ classification ဘယ်လောက်လုပ်ပေးနိုင်သလဲ ဆိုတာကို experiment လုပ်ကြည့်တာမျိုး
- အဲဒီ ရလဒ်တွေကိုမှ ငါတို့ရဲ့ proposed formula နဲ့ confirm လုပ်ကြည့်တာမျိုး
- လုအပ်ရင် formula ကို ပြန် revise လုပ်တာမျိုး (လုအပ်မယ်လို့ နားလည်)

# Project Details & Guidance (Humor Team)

- မသိသေးတဲ့ ကျောင်းသားတွေလည်း ရိုနိုင်လို့။ Literature review အတွက်သုံးပါ။

## perplexity

I am an NLP/AI researcher from Myanmar. Currently, I am planning to develop a new proposal on humor classification for the Myanmar language. First, I would like your assistance in conducting a literature review on humor detection and classification in the NLP/AI fields, covering not only English but also other languages. Please include relevant academic references and any useful code repositories.



# Project Details & Guidance (Humor Team)

- Check conda environment

```
(base) ye@lst-gpu-server-197:~/intern3/humor$ conda info --envs
# conda environments:
#
base                  * /home/ye/anaconda3
wordwranglers          /home/ye/anaconda3/envs/wordwranglers
bi_lstm_ner            /home/ye/anaconda3/envs/bi_lstm_ner
btg-seq2seq             /home/ye/anaconda3/envs/btg-seq2seq
experiments-tc          /home/ye/anaconda3/envs/experiments-tc
hs-fasttext              /home/ye/anaconda3/envs/hs-fasttext
mcrf                   /home/ye/anaconda3/envs/mcrcf
py2.7                   /home/ye/anaconda3/envs/py2.7
py3.6                   /home/ye/anaconda3/envs/py3.6
py3.7                   /home/ye/anaconda3/envs/py3.7
py3.8                   /home/ye/anaconda3/envs/py3.8
st                      /home/ye/anaconda3/envs/st
textgen                 /home/ye/anaconda3/envs/textgen

(base) ye@lst-gpu-server-197:~/intern3/humor$
```

# Project Details & Guidance (Humor Team)

- Activate an existing environment

```
(base) ye@1st-gpu-server-197:~/intern3/humor$ conda activate py3.8
(py3.8) ye@1st-gpu-server-197:~/intern3/humor$ _
```

- Create a new conda environment

```
(py3.8) ye@1st-gpu-server-197:~/intern3/humor$ conda create -n "humor" python=3.13.5
Collecting package metadata (current_repodata.json): / _
```

# Project Details & Guidance (Humor Team)

- တချို့ ဆာဗာတွေမှက ပေးမလုပ်တာမျိုးလည်း ရှိနှင့်

```
(py3.8) [ye@lst-gpu-server-197:~/intern3/humor]$ conda create -n "humor" python=3.13.5
Collecting package metadata (current_repodata.json): failed
CondaHTTPError: HTTP 000 CONNECTION FAILED for url <https://repo.anaconda.com/pkgs/main/linux-64/current_repodata.json>
Elapsed: -

An HTTP error occurred when trying to retrieve this URL.
HTTP errors are often intermittent, and a simple retry will get you on your way.

If your current network has https://www.anaconda.com blocked, please file
a support request with your network engineering team.

'https://repo.anaconda.com/pkgs/main/linux-64'
```

# Project Details & Guidance (Humor Team)

- ရုံးက GPU ဆာမာတစ်လုံးမှာ ပေးမလုပ်လို့ တခြား ဆာမာမှာ ပြောင်းလုပ်ခဲ့တယ်

```
libzlib           conda-forge/linux-64::libzlib-1.3.1-hb9d3cd8_2
ncurses          conda-forge/linux-64::ncurses-6.5-h2d0b736_3
openssl          conda-forge/linux-64::openssl-3.5.0-h7b32b05_1
pip              conda-forge/noarch::pip-25.1.1-pyh145f28c_0
python            conda-forge/linux-64::python-3.13.5-hf636f53_101_cp313
python_abi        conda-forge/noarch::python_abi-3.13.7_cp313
readline          conda-forge/linux-64::readline-8.2-h8c095d6_2
tk                conda-forge/linux-64::tk-8.6.13-noxft_hd72426e_102
tzdata           conda-forge/noarch::tzdata-2025b-h78e105d_0

Proceed ([y]/n)? y

Downloading and Extracting Packages:

Preparing transaction: done
Verifying transaction: done
Executing transaction: done
#
# To activate this environment, use
#
#     $ conda activate humor
#
# To deactivate an active environment, use
#
#     $ conda deactivate

(base) ye@lst-hpc3090:~$ conda activate humor
(humor) ye@lst-hpc3090:~$
```

# Project Details & Guidance (Humor Team)

- Kaggle binary classification humor dataset የ download ስርቃታዊ በርሃን
- pip install kaggle

```
Collecting urllib3>=1.15.1 (from kaggle)
  Downloading urllib3-2.4.0-py3-none-any.whl.metadata (6.5 kB)
Collecting webencodings (from kaggle)
  Using cached webencodings-0.5.1-py2.py3-none-any.whl.metadata (2.1 kB)
  Downloading kaggle-1.7.4.5-py3-none-any.whl (181 kB)
  Downloading certifi-2025.4.26-py3-none-any.whl (159 kB)
  Using cached python_dateutil-2.9.0.post0-py2.py3-none-any.whl (229 kB)
  Downloading setuptools-80.9.0-py3-none-any.whl (1.2 MB)
    ━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 1.2/1.2 MB 23.1 MB/s eta 0:00:00
  Using cached six-1.17.0-py2.py3-none-any.whl (11 kB)
  Downloading urllib3-2.4.0-py3-none-any.whl (128 kB)
  Using cached bleach-6.2.0-py3-none-any.whl (163 kB)
  Downloading charset_normalizer-3.4.2-cp313-cp313-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (148 kB)
  Using cached idna-3.10-py3-none-any.whl (70 kB)
  Downloading protobuf-6.31.1-cp39-abi3-manylinux2014_x86_64.whl (321 kB)
  Using cached python_slugify-8.0.4-py2.py3-none-any.whl (10 kB)
  Using cached text_unidecode-1.3-py2.py3-none-any.whl (78 kB)
  Downloading requests-2.32.4-py3-none-any.whl (64 kB)
  Using cached tqdm-4.67.1-py3-none-any.whl (78 kB)
  Using cached webencodings-0.5.1-py2.py3-none-any.whl (11 kB)
Installing collected packages: webencodings, text-unidecode, urllib3, tqdm, six, setuptools, python-slugify, protobuf, idna, charset-normalizer, certifi, bleach, requests, python-dateutil, kaggle
Successfully installed bleach-6.2.0 certifi-2025.4.26 charset-normalizer-3.4.2 idna-3.10 kaggle-1.7.4.5 protobuf-6.31.1 python-dateutil-2.9.0.post0 python-slugify-8.0.4 requests-2.32.4 setuptools-80.9.0 six-1.17.0 text-unidecode-1.3 tqdm-4.67.1 urllib3-2.4.0 webencodings-0.5.1
(humor) ye@1st-hpc3090:~$
```

# Project Details & Guidance (Humor Team)

- Kaggle account profile ဖိုင်ကို တွေ့ဥုံး user ထွေ access လုပ်လိုမရအောင်ပြောင်း

```
(humor) ye@1st-hpc3090:~$ ls .kaggle/kaggle.json  
.kaggle/kaggle.json  
(humor) ye@1st-hpc3090:~$  
(humor) ye@1st-hpc3090:~$ chmod 600 .kaggle/kaggle.json  
(humor) ye@1st-hpc3090:~$
```

# Project Details & Guidance (Humor Team)

- kaggle command သုံးလို့ ရမရ စမ်းကြည့်

category	reward	teamCount	userHasEntered	deadline	cat
<a href="https://www.kaggle.com/competitions/arc-prize-2025">https://www.kaggle.com/competitions/arc-prize-2025</a>	1,000,000 Usd	507	False	2025-11-03 23:59:00	Fee
<a href="https://www.kaggle.com/competitions/openai-to-z-challenge">https://www.kaggle.com/competitions/openai-to-z-challenge</a>	400,000 Usd	0	False	2025-06-29 23:59:00	Fee
<a href="https://www.kaggle.com/competitions/make-data-count-finding-data-references">https://www.kaggle.com/competitions/make-data-count-finding-data-references</a>	100,000 Usd	103	False	2025-09-09 23:59:00	Res
<a href="https://www.kaggle.com/competitions/waveform-inversion">https://www.kaggle.com/competitions/waveform-inversion</a>	50,000 Usd	1226	False	2025-06-30 23:59:00	Res
<a href="https://www.kaggle.com/competitions/cmi-detect-behavior-with-sensor-data">https://www.kaggle.com/competitions/cmi-detect-behavior-with-sensor-data</a>	50,000 Usd	944	False	2025-09-02 23:59:00	Fee
<a href="https://www.kaggle.com/competitions/moto-kaggle-hackathon">https://www.kaggle.com/competitions/moto-kaggle-hackathon</a>				2025-07-31 23:59:00	Fee

- သို့သော Kaggle account လုပ်ထားရမယ်
- ဒေတာ download လုပ်ဖို့ဆုံးရင် competition အတွက် agreement လိုအပ်တယ်

# Project Details & Guidance (Humor Team)

- Humor detection dataset ကို download လုပ်ခဲ့ (binary classification dataset)

The screenshot shows the Kaggle interface for a dataset titled "Humor Detection". The left sidebar includes links for Create, Home, Competitions, Datasets, Models, Code, Discussions, Learn, and More. Under "Your Work", there are entries for "Text Humor Detectio...", "Mission Torch 2 (R...)", and "learning pytorch 3: c...". Under "Viewed", there is a link to "View Active Events".

The main content area displays the "processed\_data(20k).csv" file (916.9 kB). It has three tabs: Detail (selected), Compact, and Column. The "About this file" section notes that the file does not have a description yet. A "Suggest Edits" button is present. Below this, a table shows the distribution of values in the "text" and "# humour" columns. The "# humour" column has two categories: 0 and 1.

text	# humour
19905 unique values	0 1
joe biden rule bid guy not run	0
watch darvish gave hitter whiplash slow pitch	0
call turtl without shell dead	1

The "Data Explorer" sidebar shows "Version 1 (3.21 MB)" containing "processed\_data(20k).csv" and "processed\_data(50k).csv". The "Summary" sidebar indicates there are 2 files and 4 columns.

# Project Details & Guidance (Humor Team)

```
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle$ wc *.csv
 20001 142032 916899 processed_data(20k).csv
 50001 354846 2290501 processed_data(50k).csv
 70002 496878 3207400 total
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle$ head ./processed_data\(20k\).csv
text, humour
joe biden rule bid guy not run,0
watch darvish gave hitter whiplash slow pitch,0
call turtl without shell dead,1
reason elect feel person,0
pasco polic shot mexican migrant behind new autopsi show,0
martha stewart tweet hideou food photo twitter respond accordingli,0
pokemon master favorit kind pasta wartortellini,1
nativ american hate rain april bring mayflow,1
obama climat chang legaci impress imperfect vulner,0
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle$ head ./processed_data\(50k\).csv
text, humour
joe biden rule bid guy not run,0
watch darvish gave hitter whiplash slow pitch,0
call turtl without shell dead,1
reason elect feel person,0
pasco polic shot mexican migrant behind new autopsi show,0
martha stewart tweet hideou food photo twitter respond accordingli,0
pokemon master favorit kind pasta wartortellini,1
nativ american hate rain april bring mayflow,1
obama climat chang legaci impress imperfect vulner,0
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle$ -
```

# Project Details & Guidance (Humor Team)

- 50k က 20k ကို extend လုပ်ထားတဲ့ပုံရှိတယ်
- 20k ထဲက စာကြောင်းတွေက 50k ဒေတာထဲမှာ ထပ်ပါနေတယ်
- အဲဒါကြောင့် 50k ကိုပဲ သုံးမယ်
- လက်ရှု ဖုံးနာမည်ကြီးက အဆင်မပြောဘူး၊ ဖုံးနာမည်ပြောင်းမယ်

```
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$ mv processed_data\50k\*.csv humor_50k.csv
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$ wc ./humor_50k.csv
 50001 354846 2290501 ./humor_50k.csv
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$ file ./humor_50k.csv
./humor_50k.csv: CSV Unicode text, UTF-8 text
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$
```

# Project Details & Guidance (Humor Team)

- Training data, Test data 90:10 គឺដោយខ្លួន

```
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$ shuf ./humor_50k.csv > ./humor_50k.shuf.csv
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$ head -n 45001 ./humor_50k.shuf.csv > train.csv
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$ tail -n 5000 ./humor_50k.shuf.csv > test.csv
```

```
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$ head ./train.csv
truli headscratch moment donald trump press confer infrastructur,0
indian chief friend got indian chief tattoo arm arm never work,1
destin wed tip bridesmaid budget,0
dark alley johnni optimist beat half life,1
horni pirat worst nightmar sunken chest booti,1
mathematician work home function domain,1
bought drug shoe dealer unlac still got high heel,1
virginia museum open costa concordia exhibit memori day weekend photo,0
peopl nt realiz chickpea get everi manpea make,1
betabrand hire graduat student model result pretti great,0
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$ 
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$ tail ./test.csv
rob gray man bun ncaa tournameant,0
texa put undu burden women choic abort clinic tell suprem court,0
alabama gop told gay sheriff candid run republican,0
nt phone sex might get hear aid,1
not lazi energi save mode,1
q nt blond like butter toast ca nt figur side butter goe,1
like coffe like like tea hot splash milk,1
hope nt take joke liter pleas return later,1
syrian increasingli desper escap wartorn countri,0
herb chicken thigh roast paper bag,0
(humor) ye@lst-hpc3090:~/intern3/humor/kaggle/data$
```

# Project Details & Guidance (Humor Team)

```
(humor) ye@lst-hpc3090:~/intern3/humor$ python3.13 ./ml_humor_detection.py --help
usage: ml_humor_detection.py [-h] [--train_file TRAIN_FILE] [-t test_file TEST_FILE]
                             [--classifier {all,svm,random_forest,logistic_regression,naive_bayes,knn,decision_tree,adaboost,gradient_boosting,voting}]
                             [--svm_kernel {linear,rbf,poly,sigmoid}] [--svm_c SVM_C]
                             [--rf_n_estimators RF_N_ESTIMATORS] [--rf_max_depth RF_MAX_DEPTH]
                             [--lr_c LR_C] [--lr_penalty {l1,l2,elasticnet,none}] [--nb_alpha NB_ALPHA]
                             [--knn_n_neighbors KNN_N_NEIGHBORS] [--dt_max_depth DT_MAX_DEPTH]
                             [--ab_n_estimators AB_N_ESTIMATORS] [--ab_learning_rate AB_LEARNING_RATE]
                             [--gb_n_estimators GB_N_ESTIMATORS] [--gb_learning_rate GB_LEARNING_RATE]

Humor Detection using Traditional ML Techniques

options:
  -h, --help            show this help message and exit
  --train_file TRAIN_FILE      Path to training CSV file (default: train.csv)
  --test_file TEST_FILE       Path to testing CSV file (default: test.csv)
  --classifier {all,svm,random_forest,logistic_regression,naive_bayes,knn,decision_tree,adaboost,gradient
                 _boosting,voting}          Classifier to use (default: all)
  --svm_kernel {linear,rbf,poly,sigmoid}        Kernel type for SVM (default: linear)
  --svm_c SVM_C                Regularization parameter for SVM (default: 1.0)
  --rf_n_estimators RF_N_ESTIMATORS        Number of trees in Random Forest (default: 100)
  --rf_max_depth RF_MAX_DEPTH           Maximum depth of trees in Random Forest (default: None)
  --lr_c LR_C                  Inverse of regularization strength for Logistic Regression (default: 1.0)
  --lr_penalty {l1,l2,elasticnet,none}    Penalty norm for Logistic Regression (default: l2)
  --nb_alpha NB_ALPHA            Additive smoothing parameter for Naive Bayes (default: 1.0)
  --knn_n_neighbors KNN_N_NEIGHBORS      Number of neighbors for KNN (default: 5)
  --dt_max_depth DT_MAX_DEPTH         Maximum depth for Decision Tree (default: None)
  --ab_n_estimators AB_N_ESTIMATORS      Number of estimators for AdaBoost (default: 50)
  --ab_learning_rate AB_LEARNING_RATE    Learning rate for AdaBoost (default: 1.0)
  --gb_n_estimators GB_N_ESTIMATORS      Number of estimators for Gradient Boosting (default: 100)
  --gb_learning_rate GB_LEARNING_RATE    Learning rate for Gradient Boosting (default: 0.1)
```

- Coding for machine learning approaches

# Project Details & Guidance (Humor Team)

- လိုအပ်တဲ့ python library ထွေကြု installation လုပ်ခဲ့တယ်

```
(humor) ye@lst-hpc3090:~/intern3/humor$ pip install scikit-learn
Collecting scikit-learn
  Downloading scikit_learn-1.7.0-cp313-cp313-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (17 kB)
Collecting numpy>=1.22.0 (from scikit-learn)
  Downloading numpy-2.3.0-cp313-cp313-manylinux_2_28_x86_64.whl.metadata (62 kB)
Collecting scipy>=1.8.0 (from scikit-learn)
  Downloading scipy-1.15.3-cp313-cp313-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (61 kB)
Collecting joblib>=1.2.0 (from scikit-learn)
  Downloading joblib-1.5.1-py3-none-any.whl.metadata (5.6 kB)
Collecting threadpoolctl>=3.1.0 (from scikit-learn)
  Using cached threadpoolctl-3.6.0-py3-none-any.whl.metadata (13 kB)
  Downloading scikit_learn-1.7.0-cp313-cp313-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (12.5 MB)
    ████████████████████████████████████████████ 12.5/12.5 MB 56.7 MB/s eta 0:00:00
  Downloading joblib-1.5.1-py3-none-any.whl (307 kB)
  Downloading numpy-2.3.0-cp313-cp313-manylinux_2_28_x86_64.whl (16.6 MB)
    ████████████████████████████████████████████ 16.6/16.6 MB 68.0 MB/s eta 0:00:00
  Downloading scipy-1.15.3-cp313-cp313-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (37.3 MB)
    ████████████████████████████████████████████ 37.3/37.3 MB 74.6 MB/s eta 0:00:00
  Using cached threadpoolctl-3.6.0-py3-none-any.whl (18 kB)
  Installing collected packages: threadpoolctl, numpy, joblib, scipy, scikit-learn
  Successfully installed joblib-1.5.1 numpy-2.3.0 scikit-learn-1.7.0 scipy-1.15.3 threadpoolctl-3.6.0
(humor) ye@lst-hpc3090:~/intern3/humor$ -
```

# Project Details & Guidance (Humor Team)

- Pandas installation

```
(humor) ye@lst-hpc3090:~/intern3/humor$ pip install pandas
Collecting pandas
  Downloading pandas-2.3.0-cp313-cp313-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (91 kB)
Requirement already satisfied: numpy>=1.26.0 in /home/ye/miniforge3/envs/humor/lib/python3.13/site-packages (from pandas) (2.3.0)
Requirement already satisfied: python-dateutil>=2.8.2 in /home/ye/miniforge3/envs/humor/lib/python3.13/site-packages (from pandas) (2.9.0.post0)
Collecting pytz>=2020.1 (from pandas)
  Using cached pytz-2025.2-py2.py3-none-any.whl.metadata (22 kB)
Collecting tzdata>=2022.7 (from pandas)
  Using cached tzdata-2025.2-py2.py3-none-any.whl.metadata (1.4 kB)
Requirement already satisfied: six>=1.5 in /home/ye/miniforge3/envs/humor/lib/python3.13/site-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)
Downloading pandas-2.3.0-cp313-cp313-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (12.0 MB)
  ████████████████████████████████████████████ 12.0/12.0 MB 57.8 MB/s eta 0:00:00
Using cached pytz-2025.2-py2.py3-none-any.whl (509 kB)
Using cached tzdata-2025.2-py2.py3-none-any.whl (347 kB)
Installing collected packages: pytz, tzdata, pandas
Successfully installed pandas-2.3.0 pytz-2025.2 tzdata-2025.2
(humor) ye@lst-hpc3090:~/intern3/humor$ _
```

# Project Details & Guidance (Humor Team)

- tabulate installation

```
(humor) ye@1st-hpc3090:~/intern3/humor$ pip install tabulate
Collecting tabulate
  Using cached tabulate-0.9.0-py3-none-any.whl.metadata (34 kB)
Using cached tabulate-0.9.0-py3-none-any.whl (35 kB)
Installing collected packages: tabulate
Successfully installed tabulate-0.9.0
(humor) ye@1st-hpc3090:~/intern3/humor$ -
```

# Project Details & Guidance (Humor Team)

- Training/Testing with ML approaches

Comparison of All Classifiers:					
Model	Accuracy	Precision	Recall	F1 Score	Time
Svm	0.8618	0.8619	0.8677	0.8648	243.74s
Random Forest	0.8416	0.8295	0.8673	0.848	29.79s
Logistic Regression	0.8652	0.866	0.87	0.868	0.04s
Naive Bayes	0.8658	0.8658	0.8716	0.8687	0.01s
Knn	0.6062	0.5691	0.934	0.7073	1.56s
Decision Tree	0.7998	0.8024	0.8053	0.8038	5.79s
Adaboost	0.6028	0.5639	0.9717	0.7137	1.22s
Gradient Boosting	0.7214	0.8359	0.5638	0.6734	5.97s
Voting	0.8728	0.8691	0.8834	0.8762	277.82s

Best model: Voting with F1 score: 0.8762

```
real    9m26.984s
user    9m33.239s
sys     0m1.215s
(humor) ye@1st-hpc3090:~/intern3/humor$ -
```

# Project Details & Guidance (Humor Team)

- ဆရာ လုပ်ပြခဲ့တာကို အခြေခံပြီး ML approaches တွေကို နားလည်အောင်ကြိုးစားပါ
- Neural Network based approach တွေဖြစ်တဲ့ Bi-LSTM, textCNN, Transformer တွေနဲ့ လုပ်စမ်းကြည့်ပါ
- မြန်မာစာအတွက်က ဒေတာက မရှိသေးတာမို့လို့ myHumor corpus ကို အနုည်းဆုံး binary တော့ ဆောက်ကြရှုလိမ့်မယ်
- ဖြစ်နှင့်ရှင် multi class model အတွက် class တွေကို ဘယ်လို့ သတ်မှတ်ကြရှင် ကောင်းမလဲ စဉ်းစားကြ၊ တိုင်ပင်ကြရှုလိမ့်မယ်
- ဒေတာကို ဘယ်လို့ စုကြမလဲ၊ လက်တွေ့ကြို့ရမယ့် အခက်အခဲတွေကို ဘယ်လို့ engineering နည်းလမ်းတွေနဲ့ ကျော်လွှားကြမလဲ တိုင်ပင်ကြမှ ဖြစ်လိမ့်မယ်

# Project Details & Guidance (Image Team)

- မြန်မာ အစားအသေက်ပုံတွေကို နေ့စဉ်စာတ်ပုံရှိက် စုသွားလိုလည်း ရတယ်



# Project Details & Guidance (Image Team)

- ၃ပမာ ပုဂ္ဂန်ကွက်ကြောသုပ်ဆိုရင်လည်း တစ်ဆိုင်တည်းမှာပဲ အမျိုးမျိုးရှက်နိုင်



# Project Details & Guidance (Image Team)

- မြန်မာမှန်တွေ ဆိုရင်လည်း ဗန်းထဲမှာ ထည့်ထားတဲ့ပုံ



# Project Details & Guidance (Image Team)

- ဖော်ပူးထဲမှာ ထည့်ထားတဲ့ပုံ၊ ပန်းကန်ထဲမှာ ထည့်ထားတဲ့ပုံ စသည်ဖြင့်



# Project Details & Guidance (Image Team)

- ဟင်းတွေ ဆိုရင်လည်း အိုးထဲမှာ ထည့်ထားတဲ့ပုံ၊ ချက်လက်စ အငွေ့ထွေက်နေတဲ့ပုံ



# Project Details & Guidance (Image Team)



- ထမင်းဆိုင်တွေမှာ အကူအညီတောင်းပြီး ရုက်တာမျိုးလည်း ရလိမ့်မယ်
- တရက်တည်းမှာတင် ကိုယ်စားတဲ့ မန်က်စာ၊ နေ့လည်စာ၊ ညျေနေစာ မှာ ဓာတ်ပုံစွန်း
- အရှင်ဆုံး မြန်မာအစားအသောက် ကိုပဲ ဘယ်လို အမျိုးအစား ခွဲကြမှုလဲ ဆုံးတာကို သေသေချာချာ သတ္တမှတ်ကြရလဲမ့်မယ်
- ပြီးရင် ဘာ NLP or AI task အတွက်လည်း စဉ်းစားရမယ်
- ဘယ်လို အသေးစိတ် စဉ်းစားကြရမယ်

# Project Details & Guidance (Image Team)

arXiv > cs > arXiv:2408.01723

Computer Science > Computer Vision and Pattern Recognition

[Submitted on 3 Aug 2024]

## A Novel Evaluation Framework for Image2Text Generation

REQUEST CODE   

Jia-Hong Huang, Hongyi Zhu, Yixian Shen, Stevan Rudinac, Alessio M. Pacces, Evangelos Kanoulas

Ask the author(s) a question! :)

Ask

powered by CatalyzeX

Evaluating the quality of automatically generated image descriptions is challenging, requiring metrics that capture various aspects such as grammaticality, coverage, correctness, and truthfulness. While human evaluation offers valuable insights, its cost and time-consuming nature pose limitations. Existing automated metrics like BLEU, ROUGE, METEOR, and CIDEr aim to bridge this gap but often show weak correlations with human judgment. We address this challenge by introducing a novel evaluation framework rooted in a modern large language model (LLM), such as GPT-4 or Gemini, capable of image generation. In our proposed framework, we begin by feeding an input image into a designated image captioning model, chosen for evaluation, to generate a textual description. Using this description, an LLM then creates a new image. By extracting features from both the original and LLM-created images, we measure their similarity using a designated similarity metric. A high similarity score suggests that the image captioning model has accurately generated textual descriptions, while a low similarity score indicates discrepancies, revealing potential shortcomings in the model's performance. Human-annotated reference captions are not required in our proposed evaluation framework, which serves as a valuable tool for evaluating the effectiveness of image captioning models. Its efficacy is confirmed through human evaluation.

# Next Steps & Action Items

- ပြင်ထားတဲ့ ဆလိုက်တွေကို ကြည့်ရင်ပဲ ဆရာ တတ်နိုင်သလောက် အချိန်ပေး ပြင်ဆင်ထားတယ် ဆုတာကို ခန့်မှန်းနိုင်ပါလိမ့်မယ်
- နောက် တစ်ပတ် 3rd Seminar မှာတော့ တစ်ဖွံ့ချုပ်းစီက ကိုယ့် proposal/project အတွက် ဘယ်လိုစဉ်းစားထားတယ်၊ ဘယ်လိုပုံစံနဲ့ အေတာကို စုကြမယ်၊ ဘယ်လို schedule နဲ့ သွားမယ် ဆုတာကို ဆရာ နားထောင်ချုပ်တယ်
- Project ခေါင်းစဉ် တစ်ခုချုပ်းစီအတွက် လက်ရှိ တခြားသူတွေ ဘယ်လို နည်းလမ်းတွေနဲ့ လုပ်နေကြတယ် ဆုတာကိုလည်း literature review စလုပ်ကြရအောင်
- အေတာကုလည်း စပြင်ဆင်နိုင်တဲ့သူက ပြင်ဆင်ကြရအောင်
- Senior member တွေကလည်း team member တွေနဲ့ စတင် စကားပြောတာ လုပ်ပါ

# Open Discussion & Q&A

# References

1. Dale-chall-3000-words.txt:  
<https://gist.github.com/Abhishek-P/e00edcc6f508640fe24f263f5836adc>
2. The Flesch reading ease score: why and how to use it:  
<https://yoast.com/flesch-reading-ease-score/>
3. Flesch–Kincaid readability tests:  
[https://en.wikipedia.org/wiki/Flesch%20%93Kincaid\\_readability\\_tests](https://en.wikipedia.org/wiki/Flesch%20%93Kincaid_readability_tests)
4. Yoast SEO 19.3: Schema improvements, new word complexity assessment: <https://yoast.com/yoast-seo-july-12-2022/>
5. Derivation of new readability formulas, Kincaid, J.P.; Fishburne, R.P.; Rogers, R.L.; Chissom, B.S. (1975)  
<https://apps.dtic.mil/sti/pdfs/ADA006655.pdf>

# References

6. G. M. McClure, "Readability formulas: Useful or useless?," in IEEE Transactions on Professional Communication, vol. PC-30, no. 1, pp. 12-15, March 1987, doi: 10.1109/TPC.1987.6449109. keywords: {Readability metrics;Computers;Training;Personnel;Materials;Educational institutions},
7. Book Lies: Readability is Impossible to Measure, Gabe Habash -- July 20th, 2011,  
<https://web.archive.org/web/20140521073931/http://blogs.publishersweekly.com/blogs/PWxyz/2011/07/20/book-lies-readability-is-impossible-to-measure>
8. Spacy\_readability: <https://www.kaggle.com/code/bond005/spacy-readability>
9. Setting up Kaggle API in Linux:  
<https://medium.com/@c.venkataramanan1/setting-up-kaggle-api-in-linux-b05332cde53a>