

Team members

1. Sai Hruthik Gangapuram
2. Sujith Kumar Gajarla

Task 1: Importing packages

```
In [1]: import re
import torch
import random
import pandas as pd
import torch.nn as nn
from torch.utils.data import Dataset
from torch.utils.data import DataLoader
from torch.nn.utils.rnn import pad_sequence
from sklearn.model_selection import train_test_split
device = 'cpu'
```

Task 2: Data Loading

```
In [2]: data_frame = pd.read_excel('dataset.xlsx')
data_frame.head(6)
```

Out[2]:

	English	Hindi
0	Yale offers advanced degrees through its Gradu...	येल अपने ग्रेजुएट स्कूल ऑफ आर्ट्स एंड साइंसेज ...
1	Browse the organizations below for information...	अध्ययन के कार्यक्रमों, शैक्षणिक आवश्यकताओं और ...
2	Graduate School of Arts & Sciences.	ग्रेजुएट स्कूल ऑफ आर्ट्स एंड साइंसेज।
3	Yale's Graduate School of Arts & Sciences offe...	येल के ग्रेजुएट स्कूल ऑफ आर्ट्स एंड साइंसेज एम...
4	School of Architecture.	स्कूल ऑफ आर्किटेक्चर।
5	The Yale School of Architecture's mandate is f...	येल स्कूल ऑफ आर्किटेक्चर का जनादेश प्रत्येक छा...

Task 3: Data preprocessing

1. Word to Index
2. Index to word
3. Word counts
4. Normalizing the sentence

```
In [3]: SOS_token = 0
EOS_token = 1

class Lang:
    def __init__(self, name):
        self.name = name
        self.word2index = { "SOS": SOS_token, "EOS": EOS_token }
        self.word2count = {}
        self.index2word = { SOS_token: "SOS", EOS_token: "EOS" }
        self.n_words = 2 # Count SOS and EOS

    def addSentence(self, sentence):
        for word in sentence.split(' '):
            self.addWord(word)

    def addWord(self, word):
        if word not in self.word2index:
            self.word2index[word] = self.n_words
            self.word2count[word] = 1
            self.index2word[self.n_words] = word
            self.n_words += 1
        else:
            self.word2count[word] += 1
```

```
In [4]: def normalizeString(sentence):
        sentence = sentence.lower().strip()
        sentence = sentence.replace('\xa0', ' ')
        sentence = re.sub(r"([, . ! ?])", r" \1", sentence)
        sentence = re.sub(r"[. ! ?]+", r"", sentence)
        return sentence
data_frame['English'] = data_frame['English'].apply(lambda sentence: normalizeString(sentence))
```

```
data_frame['Hindi'] = data_frame['Hindi'].apply(lambda sentence: normalizeString(sentence))
data_frame.head(5)
```

Out[4]:

	English	Hindi
0	yale offers advanced degrees through its gradu...	येल अपने ग्रेजुएट स्कूल ऑफ आर्ट्स एंड साइंसेज ...
1	browse the organizations below for information...	अध्ययन के कार्यक्रमों , शैक्षणिक आवश्यकताओं और...
2	graduate school of arts & sciences	ग्रेजुएट स्कूल ऑफ आर्ट्स एंड साइंसेज।
3	yale's graduate school of arts & sciences offe...	येल के ग्रेजुएट स्कूल ऑफ आर्ट्स एंड साइंसेज एम...
4	school of architecture	स्कूल ऑफ आर्किटेक्चर।

In [5]:

```
def readLangs(data_frame):
    pairs = [list(lang_pair) for index, lang_pair in data_frame.iterrows()]
    input_lang = Lang('English')
    output_lang = Lang('Hindi')
    return input_lang, output_lang, pairs
```

In [6]:

```
def prepareData(data_frame):
    input_lang, output_lang, pairs = readLangs(data_frame)
    print("Read %s sentence pairs" % len(pairs))
    print("Counting words...")
    for pair in pairs:
        input_lang.addSentence(pair[0])
        output_lang.addSentence(pair[1])
    print("Counted words:")
    print(input_lang.name, input_lang.n_words)
    print(output_lang.name, output_lang.n_words)
    return input_lang, output_lang, pairs
input_lang, output_lang, pairs = prepareData(data_frame)
print(random.choice(pairs))
```

Read 129 sentence pairs
 Counting words...
 Counted words:
 English 533
 Hindi 598

['access to high quality patient-centered health care is a social right , not a privilege ', 'उच्च गुणवत्ता वाले रोगी-केंद्रित स्वास्थ्य देखभाल तक पहुंच एक सामाजिक अधिकार है , विशेषाधिकार नहीं।']

In [7]:

```
input_lang.index2word[21]# index to word example in input language
```

```
Out[7]: 'organizations'
```

```
In [8]: output_lang.index2word[89]# index to word example in output Language
```

```
Out[8]: 'ने'
```

Task 4 : Creating Custom Dataset

```
In [9]: class CustomDataset(Dataset):

    def __init__(self, df):
        self.df=df

    def __len__(self):
        return len(self.df)

    def indexesFromSentence(self, lang, sentence):
        return [lang.word2index[word] for word in sentence.split(' ')]

    def tensorFromSentence(self, lang, sentence):
        indexes = self.indexesFromSentence(lang, sentence)
        indexes.append(EOS_token)
        return torch.tensor(indexes, dtype=torch.long, device=device)

    def __getitem__(self ,idx):
        languages = self.df.iloc[idx]
        input_tensor = self.tensorFromSentence(input_lang, languages['English'])
        target_tensor = self.tensorFromSentence(output_lang, languages['Hindi'])
        return input_tensor, target_tensor, languages['English'], languages['Hindi']
```

Task 5: Splitting the dataset into training | testing| validation

```
In [10]: training_data, testing_data = train_test_split(data_frame, test_size=0.2, random_state=42)

        validation_data, testing_data = train_test_split(testing_data, test_size=0.5, random_state=42)
```

```
In [11]: train_data_set = CustomDataset(training_data)
        valid_data_set = CustomDataset(validation_data)
```

```
test_data_set = CustomDataset(testing_data)
```

```
In [12]: print('Size of Training dataset: {}'.format(train_data_set.__len__()))
        print('Size of Testing dataset: {}'.format(test_data_set.__len__()))
        print('Size of Validation dataset: {}'.format(valid_data_set.__len__()))
```

```
Size of Training dataset: 103
Size of Testing dataset: 13
Size of Validation dataset: 13
```

```
In [20]: train_data_set[50]# sample
```

```
Out[20]: (tensor([368,  78, 344, 369, 164, 366,  42, 370, 371,  18,   1]),
         tensor([420, 258, 217,  30, 195, 413,  53, 421, 251, 358, 171,   1]),
         'we have been expanding international collaborations in many areas ',
         'हम कई क्षेत्रों में अंतरराष्ट्रीय सहयोग का विस्तार कर रहे हैं।')
```

Task 6: Loading dataset into Batches

```
In [14]: def collate_fn(batch):
        batch = sorted(batch, key=lambda x: len(x[0]), reverse=True)
        input_seqs, target_seqs, input_language, out_language = zip(*batch)
        # Pad the input sequences with zeros
        padded_input = pad_sequence(input_seqs, batch_first=True)
        # Pad the target sequences with zeros
        padded_target = pad_sequence(target_seqs, batch_first=True)
        return padded_input, padded_target, input_language, out_language
```

```
In [15]: train_loader = DataLoader(train_data_set, batch_size=8, shuffle=True, collate_fn=collate_fn)

        val_loader = DataLoader(valid_data_set, batch_size=8, shuffle=True, collate_fn=collate_fn)

        test_loader = DataLoader(test_data_set, batch_size=8, shuffle=True, collate_fn=collate_fn)
```

```
In [16]: print('Total number of batches in train data loader: {}'.format(len(train_loader)))
        print('Total number of batches in test data loader: {}'.format(len(test_loader)))
        print('Total number of batches in validation data loader: {}'.format(len(val_loader)))
```

```
Total number of batches in train data loader: 13
Total number of batches in test data loader: 2
Total number of batches in validation data loader: 2
```

Task 7: Visualising 1st sample in each batch

Train data loader

```
In [17]: for batch_index, packed in enumerate(train_loader):
    input_tensors, output_tensors, input_language, out_language = packed
    print("\033[1mTraining Batch number-----> {}".format(batch_index+1))
    # print the first input and output tensors along with their respective languages
    print("Input Language:", input_language[0])
    print("Input Tensor Shape:", input_tensors[0].shape)
    print("Input Tensor:", input_tensors[0])

    print("Output Language:", out_language[0])
    print("Output Tensor Shape:", output_tensors[0].shape)
    print("Output Tensor:", output_tensors[0])
    print('-----')
    print("\n")
```

Training Batch number-----> 1

Input Language: browse the organizations below for information on programs of study , academic requirements , and faculty research

Input Tensor Shape: torch.Size([19])

Input Tensor: tensor([19, 20, 21, 22, 23, 24, 25, 26, 10, 27, 28, 29, 30, 28, 14, 31, 32, 18, 1])

Output Language: अध्ययन के कार्यक्रमों , शैक्षणिक आवश्यकताओं और संकाय अनुसंधान के बारे में जानकारी के लिए नीचे दिए गए संगठनों को ब्राउज़ करें।

Output Tensor Shape: torch.Size([23])

Output Tensor: tensor([22, 14, 23, 24, 25, 26, 10, 27, 28, 14, 29, 30, 31, 14, 32, 33, 34, 35, 36, 37, 38, 39, 1])

Training Batch number-----> 2

Input Language: the yale school of architecture's mandate is for each student to understand architecture as a creative , productive , innovative , and responsible practice

Input Tensor Shape: torch.Size([26])

Input Tensor: tensor([20, 2, 9, 10, 46, 47, 48, 23, 49, 50, 35, 51, 45, 52, 37, 53, 28, 54, 28, 55, 28, 14, 56, 57, 18, 1])

Output Language: येल स्कूल ऑफ आर्किटेक्चर का जनादेश प्रत्येक छात्र के लिए एक रचनात्मक , उत्पादक , अभिनव और जिम्मेदार अभ्यास के रूप में वास्तुकला को समझने के लिए है।

Output Tensor Shape: torch.Size([29])

Output Tensor: tensor([2, 5, 6, 52, 53, 54, 55, 56, 14, 32, 57, 58, 24, 59, 24, 60, 10, 61, 62, 14, 63, 30, 64, 37, 65, 14, 32, 21, 1])

Training Batch number-----> 3

Input Language: search this site to discover the range of yale's international centers and initiatives , study abroad and exchange programs , collections , and galleries

Input Tensor Shape: torch.Size([26])

Input Tensor: tensor([196, 197, 198, 35, 199, 20, 200, 10, 33, 164, 192, 14, 201, 28, 27, 202, 14, 203, 26, 28, 204, 28, 14, 205, 18, 1])

Output Language: येल के अंतरराष्ट्रीय केंद्रों और पहलों की श्रेणी , विदेश में अध्ययन और विनिमय कार्यक्रमों , संग्रहों और दीर्घाओं की खोज के लिए इस साइट को खोजें।

Output Tensor Shape: torch.Size([28])

Output Tensor: tensor([2, 14, 195, 220, 10, 222, 129, 223, 24, 224, 30, 22, 10, 225, 23, 24, 226, 10, 227, 129, 228, 14, 32, 229, 230, 37, 231, 1])

Training Batch number-----> 4

Input Language: yale center for british art to the peabody museum of natural history and numerous smaller collections , are integral parts of teaching and open to the public

Input Tensor Shape: torch.Size([29])

Input Tensor: tensor([2, 240, 23, 241, 58, 35, 20, 242, 243, 10, 244, 62, 14, 230, 245, 204, 28, 144, 246, 247, 10, 186, 14, 248, 35, 20, 139, 18, 1])

Output Language: प्राकृतिक इतिहास के पीबोडी संग्रहालय के लिए येल सेंटर फॉर ब्रिटिश आर्ट और कई छोटे संग्रह , शिक्षण के अभिन्न अंग हैं और जनता के लिए खुले हैं।

Output Tensor Shape: torch.Size([29])

Output Tensor: tensor([273, 74, 14, 274, 275, 14, 32, 2, 276, 277, 278, 67, 10, 258, 279, 280, 24, 213, 14, 281, 282, 283, 10, 284, 14, 32, 285, 171, 1])

Training Batch number-----> 5

Input Language: yale is home to a diverse student body , with students from all 50 u s states and over 120 countries

Input Tensor Shape: torch.Size([24])

Input Tensor: tensor([2, 48, 483, 35, 37, 252, 50, 250, 28, 211, 142, 238, 194, 498, 308, 38, 18, 339, 14, 227, 499, 268, 18, 1])

Output Language: येल एक विविध छात्र निकाय का घर है , जिसमें सभी 50 अमेरिकी राज्यों और 120 से अधिक देशों के छात्र हैं।

Output Tensor Shape: torch.Size([23])

Output Tensor: tensor([2, 57, 289, 56, 287, 53, 542, 106, 24, 355, 216, 564, 529, 565, 10, 566, 16, 256, 301, 14, 56, 171, 1])

Training Batch number-----> 6

Input Language: the yale school of art has a long and distinguished history of training artists of the highest caliber

Input Tensor Shape: torch.Size([20])

Input Tensor: tensor([20, 2, 9, 10, 58, 59, 37, 60, 14, 61, 62, 10, 63, 64, 10, 20, 65, 66,

18, 1])

Output Language: येल स्कूल ऑफ आर्ट में उच्चतम क्षमता के प्रशिक्षण कलाकारों का एक लंबा और विशिष्ट इतिहास है।

Output Tensor Shape: torch.Size([20])

Output Tensor: tensor([2, 5, 6, 67, 30, 68, 69, 14, 70, 71, 53, 57, 72, 10, 73, 74, 21, 1, 0, 0])

Training Batch number-----> 7

Input Language: yale's graduate school of arts & sciences offers programs leading to m a , m s , m phil , and ph d degrees in 73 departments and programs

Input Tensor Shape: torch.Size([35])

Input Tensor: tensor([33, 8, 9, 10, 11, 12, 13, 3, 26, 34, 35, 36, 37, 18, 28, 36, 38, 18, 28, 36, 39, 18, 28, 14, 40, 41, 18, 5, 42, 43, 44, 14, 26, 18, 1])

Output Language: येल के ग्रेजुएट स्कूल ऑफ आर्ट्स एंड साइंसेज एमए , एमएस , एम फिल , और पीएचडी के लिए अग्रणी कार्यक्रम प्रदान करता है। 73 विभागों और कार्यक्रमों में डिग्री।

Output Tensor Shape: torch.Size([33])

Output Tensor: tensor([2, 14, 4, 5, 6, 7, 8, 9, 41, 24, 42, 24, 43, 44, 24, 10, 45, 14, 32, 46, 47, 19, 20, 21, 48, 49, 10, 23, 30, 50, 1, 0, 0])

Training Batch number-----> 8

Input Language: yale is a member of the ivy league , a group of eight prestigious universities in the northeastern united states

Input Tensor Shape: torch.Size([22])

Input Tensor: tensor([2, 48, 37, 462, 10, 20, 463, 464, 28, 37, 465, 10, 466, 455, 316, 42, 20, 467, 338, 339, 18, 1])

Output Language: येल पूर्वोत्तर संयुक्त राज्य अमेरिका में आठ प्रतिष्ठित विश्वविद्यालयों के समूह आइवी लीग का सदस्य है।

Output Tensor Shape: torch.Size([17])

Output Tensor: tensor([2, 522, 382, 383, 384, 30, 523, 512, 359, 14, 524, 525, 526, 53, 527, 21, 1])

Training Batch number-----> 9

Input Language: yale's international research , teaching , and learning activities are undertaken in a wide variety of centers and programs across all academic fields

Input Tensor Shape: torch.Size([25])

Input Tensor: tensor([33, 164, 32, 28, 186, 28, 14, 187, 188, 144, 189, 42, 37, 190, 191, 10, 192, 14, 26, 193, 194, 29, 195, 18, 1])

Output Language: येल के अंतरराष्ट्रीय अनुसंधान , शिक्षण और सीखने की गतिविधियां सभी शैक्षणिक क्षेत्रों में विभिन्न प्रकार के केंद्रों और कार्यक्रमों में की जाती हैं।

Output Tensor Shape: torch.Size([25])

Output Tensor: tensor([2, 14, 195, 28, 24, 213, 10, 214, 129, 215, 216, 25, 217, 30, 218, 219, 14, 220, 10, 23, 30, 129, 221, 171, 1])

Training Batch number-----> 10

Input Language: yale is known for its residential college system , which provides students with a supportive community and numerous opportunities for social and intellectual engagement

Input Tensor Shape: torch.Size([26])

Input Tensor: tensor([2, 48, 460, 23, 7, 506, 277, 507, 28, 254, 508, 142, 211, 37, 509, 171, 14, 230, 137, 23, 178, 14, 510, 377, 18, 1])

Output Language: येल अपनी आवासीय कॉलेज प्रणाली के लिए जाना जाता है , जो छात्रों को एक सहायक समुदाय और सामाजिक और बौद्धिक जुड़ाव के कई अवसर प्रदान करता है।

Output Tensor Shape: torch.Size([29])

Output Tensor: tensor([2, 324, 572, 316, 573, 14, 32, 521, 245, 106, 24, 102, 261, 37, 57, 574, 308, 10, 114, 10, 575, 427, 14, 258, 163, 19, 20, 21, 1])

Training Batch number-----> 11

Input Language: the world in every theatrical discipline , creating bold art that engages the mind and delights the senses

Input Tensor Shape: torch.Size([20])

Input Tensor: tensor([20, 82, 42, 83, 84, 85, 28, 86, 87, 58, 88, 89, 20, 90, 14, 91, 20, 92, 18, 1])

Output Language: हर नाट्य विधा में दुनिया , साहसिक कला का निर्माण जो मन को आकर्षित करती है और इंद्रियों को प्रसन्न करती है।

Output Tensor Shape: torch.Size([23])

Output Tensor: tensor([95, 96, 97, 30, 98, 24, 99, 100, 53, 101, 102, 103, 37, 104, 105, 106, 10, 107, 37, 108, 105, 21, 1])

Training Batch number-----> 12

Input Language: the school of the environment is dedicated to sustaining and restoring the long-term health of the biosphere and the well-being of its people

Input Tensor Shape: torch.Size([25])

Input Tensor: tensor([20, 9, 10, 20, 105, 48, 106, 35, 107, 14, 108, 20, 109, 110, 10, 20, 111, 14, 20, 112, 10, 7, 113, 18, 1])

Output Language: पर्यावरण विद्यालय जीवमंडल के दीर्घकालिक स्वास्थ्य और इसके लोगों की भलाई को बनाए रखने और बहाल करने के लिए समर्पित है।

Output Tensor Shape: torch.Size([27])

Output Tensor: tensor([122, 123, 124, 14, 125, 126, 10, 127, 128, 129, 130, 37, 131, 132, 10, 133, 117, 14, 32, 134, 21, 1, 0, 0, 0, 0, 0])

Training Batch number-----> 13

Input Language: students , scholars , and faculty have access to over 15 million volumes as well as digital databases , and a variety of research tools

Input Tensor Shape: torch.Size([27])

Input Tensor: tensor([142, 28, 69, 28, 14, 31, 78, 174, 35, 227, 231, 232, 233, 52, 214, 52, 234, 235, 28, 14, 37, 191, 10, 32, 236, 18, 1])

Output Language: छात्रों , विद्वानों और फैकल्टी के पास 15 मिलियन से अधिक संस्करणों के साथ-साथ डिजिटल डेटाबेस और विभिन्न प्रकार के शोध उपकरण हैं।

```
Output Tensor Shape: torch.Size([24])
Output Tensor: tensor([261, 24, 79, 10, 166, 14, 262, 263, 264, 16, 256, 265, 14, 241,
266, 267, 10, 218, 219, 14, 268, 269, 171, 1])
```

Test data loader

```
In [18]: for batch_index, packed in enumerate(test_loader):
        input_tensors, output_tensors, input_language, out_language = packed
        print("\033[1mTesting Batch number-----> {}".format(batch_index+1))
        # print the first input and output tensors along with their respective languages
        print("Input Language:", input_language[0])
        print("Input Tensor Shape:", input_tensors[0].shape)
        print("Input Tensor:", input_tensors[0])

        print("Output Language:", out_language[0])
        print("Output Tensor Shape:", output_tensors[0].shape)
        print("Output Tensor:", output_tensors[0])
        print('-----')
        print("\n")
```

Testing Batch number-----> 1

Input Language: opportunities for study or research abroad as well as exchange programs are managed by the individual schools and programs

Input Tensor Shape: torch.Size([21])

Input Tensor: tensor([137, 23, 27, 213, 32, 202, 52, 214, 52, 203, 26, 144, 215, 216, 20, 217, 17, 14, 26, 18, 1])

Output Language: विदेशों में अध्ययन या अनुसंधान के अवसरों के साथ-साथ विनिमय कार्यक्रमों का प्रबंधन व्यक्तिगत स्कूलों और कार्यक्रमों द्वारा किया जाता है।

Output Tensor Shape: torch.Size([22])

Output Tensor: tensor([238, 30, 22, 239, 28, 14, 240, 14, 241, 225, 23, 53, 164, 242, 13, 10, 23, 243, 244, 245, 21, 1])

Testing Batch number-----> 2

Input Language: yale offers significant financial assistance to international students to cover tuition costs as it does with students from the u s

Input Tensor Shape: torch.Size([23])

Input Tensor: tensor([2, 3, 302, 287, 220, 35, 164, 142, 35, 303, 304, 305, 52, 306, 307, 211, 142, 238, 20, 308, 38, 18, 1])

Output Language: येल अंतरराष्ट्रीय छात्रों को ट्यूशन की लागत को कवर करने के लिए महत्वपूर्ण वित्तीय सहायता प्रदान करता है जैसा कि यह यू एस के छात्रों के साथ करता है।

Output Tensor Shape: torch.Size([31])

```
Output Tensor: tensor([ 2, 195, 261, 37, 339, 129, 340, 37, 341, 117, 14, 32, 342, 323,
                        249, 19, 20, 106, 343, 344, 345, 346, 347, 348, 14, 261, 14, 235,
                        20, 21, 1])
```

Validation data loader

```
In [19]: for batch_index, packed in enumerate(val_loader):
        input_tensors, output_tensors, input_language, out_language = packed
        print("\033[1mValidation Batch number-----> {}\033[0m".format(batch_index+1))
        # print the first input and output tensors along with their respective languages
        print("Input Language:", input_language[0])
        print("Input Tensor Shape:", input_tensors[0].shape)
        print("Input Tensor:", input_tensors[0])

        print("Output Language:", out_language[0])
        print("Output Tensor Shape:", output_tensors[0].shape)
        print("Output Tensor:", output_tensors[0])
        print('-----')
        print("\n")
```

Validation Batch number-----> 1

Input Language: the jackson school of global affairs trains and equips a new generation of leaders to devise thoughtful , evidence-based solutions for challenging global problems

Input Tensor Shape: torch.Size([26])

Input Tensor: tensor([20, 114, 9, 10, 102, 115, 116, 14, 117, 37, 118, 119, 10, 72, 35, 120, 121, 28, 122, 123, 23, 124, 102, 104, 18, 1])

Output Language: जैक्सन स्कूल ऑफ ग्लोबल अफेयर्स चुनौतीपूर्ण वैश्विक समस्याओं के लिए विचारशील , साक्ष्य-आधारित समाधान तैयार करने के लिए नेताओं की एक नई पीढ़ी को प्रशिक्षित और सुसज्जित करता है।

Output Tensor Shape: torch.Size([30])

Output Tensor: tensor([135, 5, 6, 136, 138, 139, 113, 115, 14, 32, 140, 24, 141, 116, 142, 117, 14, 32, 82, 129, 57, 143, 144, 37, 145, 10, 146, 20, 21, 1])

Validation Batch number-----> 2

Input Language: today , yale welcomes the largest international community in its history , with a current enrollment of 2 ,841 international students from 121 countries

Input Tensor Shape: torch.Size([26])

Input Tensor: tensor([260, 28, 2, 261, 20, 262, 164, 171, 42, 7, 62, 28, 211, 37, 263, 264, 10, 265, 266, 164, 142, 238, 267, 268, 18, 1])

Output Language: आज , येल 121 देशों के 2 ,841 अंतराष्ट्रीय छात्रों के वर्तमान नामांकन के साथ , अपने इतिहास में सबसे बड़े अंतरराष्ट्रीय समुदाय का स्वागत करता

है।

Output Tensor Shape: torch.Size([28])

Output Tensor: tensor([299, 24, 2, 300, 301, 14, 302, 303, 292, 261, 14, 304, 305, 14,
235, 24, 3, 74, 30, 306, 307, 195, 308, 53, 309, 20, 21, 1])
