



PREDICTIVE TEXT FOR CONTEXT ANALYZATION

Gadicherla Sameer
01FB15ECS101
G Manasa
01FB15ECS099
UE15CS333 project
submission

SUBMISSION GUIDELINES

1. Since ESA time table is preponed (starting on 9th May), last date of submission is 6th May 2018 (Sunday) EOD. No extension will be granted.
2. Dataset, working ipython notebook and this presentation have to be submitted as a zip file (only one zip file per team). Name the zip file appropriately. You can either upload in Google drive or mail to bhaskarjyoti01@gmail.com
3. The final marks awarded by the evaluator **will be based on the following factors.**
 - a) Working code that delivers result : 12
 - b) Visualization, metric, comparison : 3
 - c) Nature of dataset : 2
 - d) Novelty of idea / publishability if completed : 3

ABOUT THE PROJECT

Project idea: To create a system which, given input as a part of a complete English sentence, can predict the next word in the sentence or generate the next sequence of words using Markov Chains and RNNs.

Why you think this can be turned into a publishable work:

This can be turned into a publishable work because analyzing the context of any text for a given situation is done unbiased without any human intervention. Usually this is done character by character but we are doing it word by word, which implies the dependency on previous words is captured than the previous character, which helps in remembering the context better because we train the RNN on the words directly. We will use this model in future to predict musical compositions of a particular singer. (This can be achieved even by training it on note by note in music as character by character)

LITERATURE SURVEY AND EXISTING WORK

Serial No	Title, author, web URL (if any)	What was useful as a reference
1	https://en.wikipedia.org/wiki/Language_model	We came to know what language models are exactly
2	https://en.wikipedia.org/wiki/Markov_property	Markovian property
3	https://en.wikipedia.org/wiki/Markov_chain	Markov chains
4	https://en.wikipedia.org/wiki/Recurrent_neural_network	RNN basics
5	Bengio, Yoshua, Simard, Patrice, and Frasconi, Paolo. Learning long-term dependencies with gradient descent is difficult. Neural Networks, IEEE Transactions on, 5(2):157–166, 1994.	Short term memory to train gradient descent
6	Dauphin, Yann N, de Vries, Harm, Chung, Junyoung, and Bengio, Yoshua. Rmsprop and equilibrated adaptive learning rates for non-convex optimization. arXiv preprint arXiv:1502.04390, 2015.	Rmsprop optimization though we used adam
7	Graves, Alex. Generating sequences with recurrent neural networks. arXiv preprint arXiv:1308.0850, 2013.	Basic idea on how to generate sequences
8	http://cs229.stanford.edu/proj2015/054_report.pdf	Adam optimization and normalizations

DATASET SOURCE AND PREPROCESSING DONE

1. Dataset source :

- For markov chains :

1. Corpus of Contemporary American English (COCA) from here

- About Dataset: The Corpus of Contemporary American English (COCA) is the largest freely-available corpus of English, and the only large and balanced corpus of American English.
- Attributes: We will be using a subset of COCA mainly from [here](#). Hence, we will be using two, three, four and five-gram inputs to train the model each time.
- Instance count: Each of the n-gram datasets contains the following number of n-grams:
 - 2-gram: 1,020,386
 - 3-gram: 1,020,010
 - 4-gram: 1,034,308
 - 5-gram: 1,044,269
- This was the primary dataset for testing the goodness of fit for every other dataset as it was exhaustive. We also trained using this dataset for next-word prediction as well as generating text.

1. Wikipedia Corpus

- About dataset: This corpus contains the full text of Wikipedia, and it contains 1.9 billion words in more than 4.4 million articles.
- Attributes: Wikipedia Corpus contains Wikipedia pages scraped and dumped into text files. Since it is only for testing, we used random articles with a total of 2286765 words.

2. RNN and LSTM:

1. Narendra Modi speech in US : This has around 1000+ lines and is taken from a news website
2. Narendra Modi speech on gst: This has around 2500+ lines and is also taken from a blog on web.
3. Harry Potter Novel : It is the first book of the Harry potter series with has around 7000 lines and is has around 57 distinct words.

DATASET SOURCE AND PREPROCESSING DONE(Cont.)

1. Pre-processing steps performed

1. Wikipedia Corpus:

- Read all the lines of the input file.
- Remove all unwanted characters and punctuations.
- Remove lines with very few words.
- Convert everything to lower-case.
- Run a sliding window on the words

2. Narendra Modi Speeches:

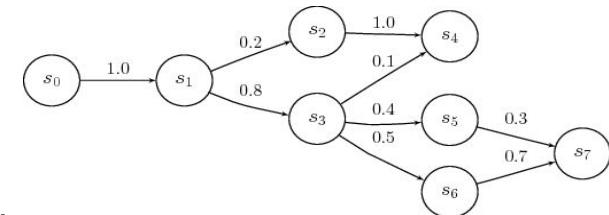
- Same as for Wikipedia but further,
 - Distinct number of words are found
 - All words are represented as numbers
 - All sentences are represented as a sequence of numbers
 - All the sequences are padded to the largest sequence size with zeroes

MARKOV CHAINS

HIGH LEVEL DESIGN OF OUR IMPLEMENTATION

MARKOV CHAIN APPROACH:

- We create the training set by creating (n-gram , n+1 word , count) pairs and save these pairs in the SQL database
- Then as said in pre-processing we also use the Wikipedia corpus, run. A sliding window and save the input,output pairs , which can also be used as training dataset but we are just using to test our model
- So , in order to predict the next word , we give a simple SQL query to the database and return an ordered dictionary of word , count pairs which the client.py file will help us calculate the probability(if needed).
- We have used the database approach just because , it is very fast compared to vanilla lists , and practically impossible when the size of the text moves into GBs.



Note : Both the dataset and latest jupyter notebook (python 3.0) have to be submitted for evaluation

PSEUDO-CODE

Create bi-grams:

```
1. >>> s = "I am Sam. Sam I am. I do not like green eggs and ham."
2. >>> tokens = s.split(" ")
3. >>> bigrams = [(tokens[i], tokens[i + 1]) for i in range(0, len(tokens) - 1)]
4. >>> bigrams[('I', 'am'), ('am', 'Sam.'), ('Sam.', 'Sam'), ('Sam', 'I'), ('I', 'am.'), ('am.', 'I'), ('I', 'do'), ('do', 'not'), ('not', 'like'), ('like', 'green'), ('green', 'eggs'), ('eggs', 'and'), ('and', 'ham.')]

```

Markov Chain: I am Sam.

```
1. {
2.     'I': ['am'],
3.     'am': ['Sam.'],
4. }

```

Add Sam I am.

```
1. {
2.     'I': ['am', 'am.'],
3.     'am': ['Sam.'],
4.     'Sam': ['I'],
5. }

```

Thus, give an input 'am', we get 'Sam.'

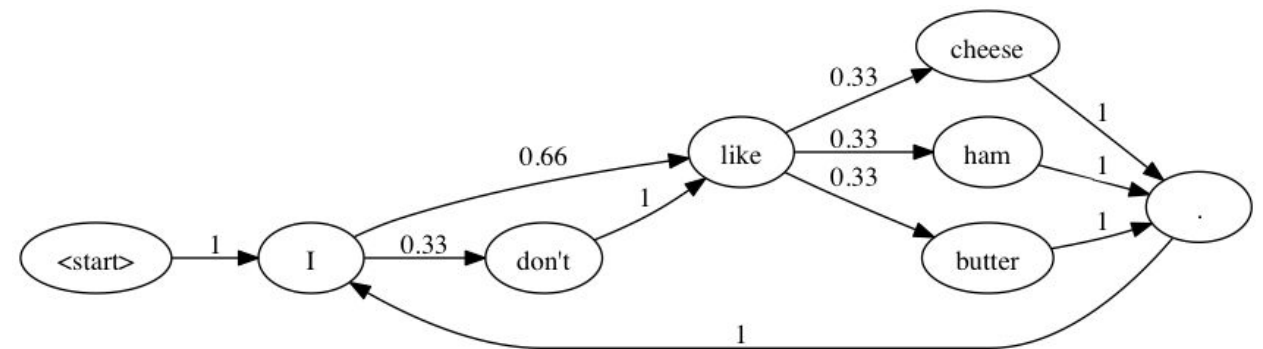


Figure 2 An example of fully modeled Markov Chain

RNN AND LSTM

RNN basics

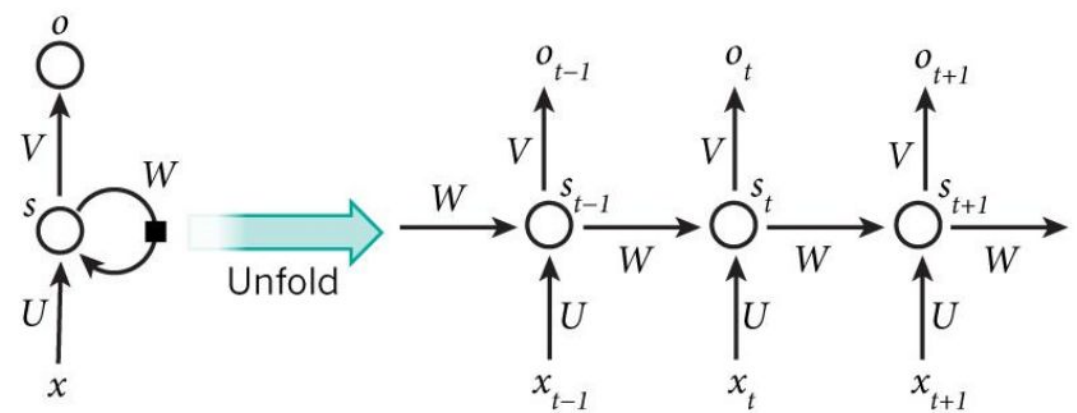


Figure 2 A typical RNN looks like one shown above. A RNN is unrolled (or unfolded) over time into a full network.

RNNs make use of sequential information. A traditional neural network assumes that all inputs (and outputs) are independent of each other.

Since our project is to predict text, it is crucial to know the information from previous states, i.e. the previous words to predict the next word.

RNNs are called recurrent because they perform the same task for every element of a sequence, with the output depending on the previous computations.

RNNs thus have a “memory” element to capture information from previous calculations.

Need for LSTM and its advantages over RNN:

Two major problems with RNN: Vanishing and Exploding gradients.

- Vanishing Gradients: In traditional RNNs the gradient signal can be multiplied a large number of times by the weight matrix. Thus, the magnitude of the weights of the transition matrix can play an important role. If the weights in the matrix are small, the gradient signal becomes smaller at every training step, thus making learning very slow or completely stops it. This is called vanishing gradient.
- Exploding Gradients: refers to the weights in this matrix being so large that it can cause learning to diverge.

LSTM stands for Long Short Term Memory, is a special kind of RNN that learns long-term dependencies. The memory cell of LSTM is composed of four elements: input, forget and output gates and a neuron that connects to itself.

CHARACTER MODEL

Design For the Character Model

In order to train our recurrent neural network, we use a large text file (we have used the Harry Potter books) as input. This large text file undergoes preprocessing and is stored in the form of two arrays X and Y defined as:

X:

which is a three dimensional array of the form $X[i, t, \text{char_indices}[\text{char}]]$ where i is the index of each sentence, t is the index of each letter in each sentence and $\text{char_indices}[\text{char}]$ indicates the letter encountered. Therefore, X is like the given input matrix.

Y:

which is a two dimensional array of the form $Y[i, \text{char_indices}[\text{next_chars}[i]]]$ where i is the index of each sentence and $\text{char_indices}[\text{next_chars}[i]]$ indicates the letter which is encountered after encountering `SEQUENCE_LENGTH` (set to 40) characters of the sentence. Therefore, Y is like the expected output matrix.

We iterate through the text file at a step size of 3 characters, scanning through all characters, till `SEQUENCE_LENGTH` (store in list of sentences) and storing the next expected character after `SEQUENCE_LENGTH` (store in list of next_char). This is then used in the creation of the X and Y matrices.

Model Architecture

We train the following model on the input format specified before:

```
model = Sequential()  
model.add(LSTM(128, input_shape=(SEQUENCE_LENGTH, len(chars))))  
model.add(Dense(len(chars)))  
model.add(Activation('softmax'))  
model.summary()
```

Layer (type)	Output Shape	Param #
lstm_2 (LSTM)	(None, 128)	86528
dense_2 (Dense)	(None, 40)	5160
activation_2 (Activation)	(None, 40)	0
Total params: 91,688		
Trainable params: 91,688		
Non-trainable params: 0		

Our model is trained for many number of epochs till 120 using RMSProp optimizer and uses 5% of the data for validation.

WORD BY WORD MODEL

LSTM APPROACH FOR WORD BY WORD

1. After all the pre-processing as mentioned in the pre-processing slide.
2. Then we create a RNN model which is actually an LSTM preserved model with following architecture and train/use this with two approaches.

```
In [12]: model = Sequential()
model.add(Embedding(vocab_size,128, input_length=max_length-1))
model.add(LSTM(200, recurrent_dropout = 0.2, dropout = 0.2))
model.add(Dense(vocab_size, activation='softmax'))
print(model.summary())
```

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, 57, 128)	64640
lstm_1 (LSTM)	(None, 200)	263200
dense_1 (Dense)	(None, 505)	101505

=====
Total params: 429,345
Trainable params: 429,345
Non-trainable params: 0
=====
None

APPROACHES

1. Whole sentence is used as X (input) and the next word is the output (y). This approach is better than just the previous word. That is how the data to train is created.

```
# split into input and output elements
sequences = array(sequences)
X, y = sequences[:, :-1], sequences[:, -1]
y = to_categorical(y, num_classes=vocab_size)
```

2. In this approach we create the dataset, by training it as X which has the previous two words and the y has the next output word.

```
sequences = list()
for i in range(2, len(encoded)):
    sequence = encoded[i-2:i+1]
    sequences.append(sequence)
print('Total Sequences: %d' % len(sequences))
print("\n", sequences)
```

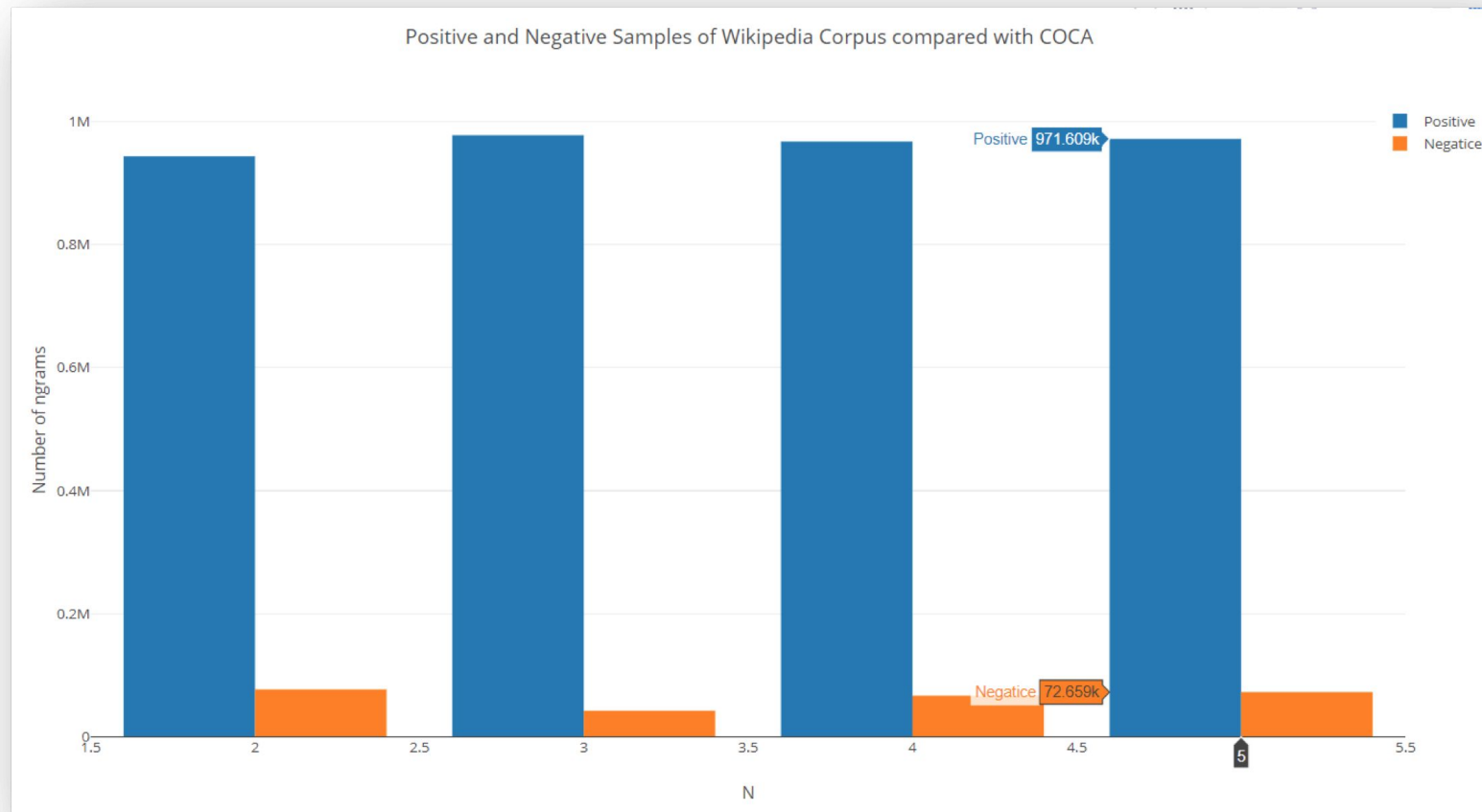
Total Sequences: 1192

```
[[148, 149, 150], [149, 150, 151], [150, 151, 152], [151, 152, 153], [152, 153, 154], [153, 154, 155], [154, 155, 156], [155, 156, 157], [156, 157, 158], [157, 158, 159], [158, 159, 160], [159, 160, 161], [160, 161, 162], [161, 162, 163], [162, 163, 164], [163, 164, 165], [164, 165, 166], [165, 166, 167], [166, 167, 168], [167, 168, 169], [168, 169, 170], [169, 170, 171], [170, 171, 172], [171, 172, 173], [172, 173, 174], [173, 174, 175], [174, 175, 176], [175, 176, 177], [176, 177, 178], [177, 178, 179], [178, 179, 180], [179, 180, 181], [180, 181, 182], [181, 182, 183], [182, 183, 184], [183, 184, 185], [184, 185, 186], [185, 186, 187], [186, 187, 188], [187, 188, 189], [188, 189, 190], [189, 190, 191], [190, 191, 192], [191, 192, 193], [192, 193, 194], [193, 194, 195], [194, 195, 196], [195, 196, 197], [196, 197, 198], [197, 198, 199], [198, 199, 200], [200, 201, 202], [201, 202, 203], [202, 203, 204], [203, 204, 205], [204, 205, 206], [205, 206, 207], [206, 207, 208], [207, 208, 209], [208, 209, 210], [209, 210, 211], [210, 211, 212], [211, 212, 213], [212, 213, 214], [213, 214, 215], [214, 215, 216], [215, 216, 217], [216, 217, 218], [217, 218, 219], [218, 219, 220], [219, 220, 221], [220, 221, 222], [221, 222, 223], [222, 223, 224], [223, 224, 225], [224, 225, 226], [225, 226, 227], [226, 227, 228], [227, 228, 229], [228, 229, 230], [229, 230, 231], [230, 231, 232], [231, 232, 233], [232, 233, 234], [233, 234, 235], [234, 235, 236], [235, 236, 237], [236, 237, 238], [237, 238, 239], [238, 239, 240], [239, 240, 241], [240, 241, 242], [241, 242, 243], [242, 243, 244], [243, 244, 245], [244, 245, 246], [245, 246, 247], [246, 247, 248], [247, 248, 249], [248, 249, 250], [249, 250, 251], [250, 251, 252], [251, 252, 253], [252, 253, 254], [253, 254, 255], [254, 255, 256], [255, 256, 257], [256, 257, 258], [257, 258, 259], [258, 259, 260], [259, 260, 261], [260, 261, 262], [261, 262, 263], [262, 263, 264], [263, 264, 265], [264, 265, 266], [265, 266, 267], [266, 267, 268], [267, 268, 269], [268, 269, 270], [269, 270, 271], [270, 271, 272], [271, 272, 273], [272, 273, 274], [273, 274, 275], [274, 275, 276], [275, 276, 277], [276, 277, 278], [277, 278, 279], [278, 279, 280], [279, 280, 281], [280, 281, 282], [281, 282, 283], [282, 283, 284], [283, 284, 285], [284, 285, 286], [285, 286, 287], [286, 287, 288], [287, 288, 289], [288, 289, 290], [289, 290, 291], [290, 291, 292], [291, 292, 293], [292, 293, 294], [293, 294, 295], [294, 295, 296], [295, 296, 297], [296, 297, 298], [297, 298, 299], [298, 299, 300], [300, 301, 302], [301, 302, 303], [302, 303, 304], [303, 304, 305], [304, 305, 306], [305, 306, 307], [306, 307, 308], [307, 308, 309], [308, 309, 310], [309, 310, 311], [310, 311, 312], [311, 312, 313], [312, 313, 314], [313, 314, 315], [314, 315, 316], [315, 316, 317], [316, 317, 318], [317, 318, 319], [318, 319, 320], [319, 320, 321], [320, 321, 322], [321, 322, 323], [322, 323, 324], [323, 324, 325], [324, 325, 326], [325, 326, 327], [326, 327, 328], [327, 328, 329], [328, 329, 330], [329, 330, 331], [330, 331, 332], [331, 332, 333], [332, 333, 334], [333, 334, 335], [334, 335, 336], [335, 336, 337], [336, 337, 338], [337, 338, 339], [338, 339, 340], [339, 340, 341], [340, 341, 342], [341, 342, 343], [342, 343, 344], [343, 344, 345], [344, 345, 346], [345, 346, 347], [346, 347, 348], [347, 348, 349], [348, 349, 350], [349, 350, 351], [350, 351, 352], [351, 352, 353], [352, 353, 354], [353, 354, 355], [354, 355, 356], [355, 356, 357], [356, 357, 358], [357, 358, 359], [358, 359, 360], [359, 360, 361], [360, 361, 362], [361, 362, 363], [362, 363, 364], [363, 364, 365], [364, 365, 366], [365, 366, 367], [366, 367, 368], [367, 368, 369], [368, 369, 370], [369, 370, 371], [370, 371, 372], [371, 372, 373], [372, 373, 374], [373, 374, 375], [374, 375, 376], [375, 376, 377], [376, 377, 378], [377, 378, 379], [378, 379, 380], [379, 380, 381], [380, 381, 382], [381, 382, 383], [382, 383, 384], [383, 384, 385], [384, 385, 386], [385, 386, 387], [386, 387, 388], [387, 388, 389], [388, 389, 390], [389, 390, 391], [390, 391, 392], [391, 392, 393], [392, 393, 394], [393, 394, 395], [394, 395, 396], [395, 396, 397], [396, 397, 398], [397, 398, 399], [398, 399, 400], [400, 401, 402], [401, 402, 403], [402, 403, 404], [403, 404, 405], [404, 405, 406], [405, 406, 407], [406, 407, 408], [407, 408, 409], [408, 409, 410], [409, 410, 411], [410, 411, 412], [411, 412, 413], [412, 413, 414], [413, 414, 415], [414, 415, 416], [415, 416, 417], [416, 417, 418], [417, 418, 419], [418, 419, 420], [419, 420, 421], [420, 421, 422], [421, 422, 423], [422, 423, 424], [423, 424, 425], [424, 425, 426], [425, 426, 427], [426, 427, 428], [427, 428, 429], [428, 429, 430], [429, 430, 431], [430, 431, 432], [431, 432, 433], [432, 433, 434], [433, 434, 435], [434, 435, 436], [435, 436, 437], [436, 437, 438], [437, 438, 439], [438, 439, 440], [439, 440, 441], [440, 441, 442], [441, 442, 443], [442, 443, 444], [443, 444, 445], [444, 445, 446], [445, 446, 447], [446, 447, 448], [447, 448, 449], [448, 449, 450], [449, 450, 451], [450, 451, 452], [451, 452, 453], [452, 453, 454], [453, 454, 455], [454, 455, 456], [455, 456, 457], [456, 457, 458], [457, 458, 459], [458, 459, 460], [459, 460, 461], [460, 461, 462], [461, 462, 463], [462, 463, 464], [463, 464, 465], [464, 465, 466], [465, 466, 467], [466, 467, 468], [467, 468, 469], [468, 469, 470], [469, 470, 471], [470, 471, 472], [471, 472, 473], [472, 473, 474], [473, 474, 475], [474, 475, 476], [475, 476, 477], [476, 477, 478], [477, 478, 479], [478, 479, 480], [479, 480, 481], [480, 481, 482], [481, 482, 483], [482, 483, 484], [483, 484, 485], [484, 485, 486], [485, 486, 487], [486, 487, 488], [487, 488, 489], [488, 489, 490], [489, 490, 491], [490, 491, 492], [491, 492, 493], [492, 493, 494], [493, 494, 495], [494, 495, 496], [495, 496, 497], [496, 497, 498], [497, 498, 499], [498, 499, 500], [500, 501, 502], [501, 502, 503], [502, 503, 504], [503, 504, 505], [504, 505, 506], [505, 506, 507], [506, 507, 508], [507, 508, 509], [508, 509, 510], [509, 510, 511], [510, 511, 512], [511, 512, 513], [512, 513, 514], [513, 514, 515], [514, 515, 516], [515, 516, 517], [516, 517, 518], [517, 518, 519], [518, 519, 520], [519, 520, 521], [520, 521, 522], [521, 522, 523], [522, 523, 524], [523, 524, 525], [524, 525, 526], [525, 526, 527], [526, 527, 528], [527, 528, 529], [528, 529, 530], [529, 530, 531], [530, 531, 532], [531, 532, 533], [532, 533, 534], [533, 534, 535], [534, 535, 536], [535, 536, 537], [536, 537, 538], [537, 538, 539], [538, 539, 540], [539, 540, 541], [540, 541, 542], [541, 542, 543], [542, 543, 544], [543, 544, 545], [544, 545, 546], [545, 546, 547], [546, 547, 548], [547, 548, 549], [548, 549, 550], [549, 550, 551], [550, 551, 552], [551, 552, 553], [552, 553, 554], [553, 554, 555], [554, 555, 556], [555, 556, 557], [556, 557, 558], [557, 558, 559], [558, 559, 560], [559, 560, 561], [560, 561, 562], [561, 562, 563], [562, 563, 564], [563, 564, 565], [564, 565, 566], [565, 566, 567], [566, 567, 568], [567, 568, 569], [568, 569, 570], [569, 570, 571], [570, 571, 572], [571, 572, 573], [572, 573, 574], [573, 574, 575], [574, 575, 576], [575, 576, 577], [576, 577, 578], [577, 578, 579], [578, 579, 580], [579, 580, 581], [580, 581, 582], [581, 582, 583], [582, 583, 584], [583, 584, 585], [584, 585, 586], [585, 586, 587], [586, 587, 588], [587, 588, 589], [588, 589, 590], [589, 590, 591], [590, 591, 592], [591, 592, 593], [592, 593, 594], [593, 594, 595], [594, 595, 596], [595, 596, 597], [596, 597, 598], [597, 598, 599], [598, 599, 600], [600, 601, 602], [601, 602, 603], [602, 603, 604], [603, 604, 605], [604, 605, 606], [605, 606, 607], [606, 607, 608], [607, 608, 609], [608, 609, 610], [609, 610, 611], [610, 611, 612], [611, 612, 613], [612, 613, 614], [613, 614, 615], [614, 615, 616], [615, 616, 617], [616, 617, 618], [617, 618, 619], [618, 619, 620], [619, 620, 621], [620, 621, 622], [621, 622, 623], [622, 623, 624], [623, 624, 625], [624, 625, 626], [625, 626, 627], [626, 627, 628], [627, 628, 629], [628, 629, 630], [629, 630, 631], [630, 631, 632], [631, 632, 633], [632, 633, 634], [633, 634, 635], [634, 635, 636], [635, 636, 637], [636, 637, 638], [637, 638, 639], [638, 639, 640], [639, 640, 641], [640, 641, 642], [641, 642, 643], [642, 643, 644], [643, 644, 645], [644, 645, 646], [645, 646, 647], [646, 647, 648], [647, 648, 649], [648, 649, 650], [649, 650, 651], [650, 651, 652], [651, 652, 653], [652, 653, 654], [653, 654, 655], [654, 655, 656], [655, 656, 657], [656, 657, 658], [657, 658, 659], [658, 659, 660], [659, 660, 661], [660, 661, 662], [661, 662, 663], [662, 663, 664], [663, 664, 665], [664, 665, 666], [665, 666, 667], [666, 667, 668], [667, 668, 669], [668, 669, 670], [669, 670, 671], [670, 671, 672], [671, 672, 673], [672, 673, 674], [673, 674, 675], [674, 675, 676], [675, 676, 677], [676, 677, 678], [677, 678, 679], [678, 679, 680], [679, 680, 681], [680, 681, 682], [681, 682, 683], [682, 683, 684], [683, 684, 685], [684, 685, 686], [685, 686, 687], [686, 687, 688], [687, 688, 689], [688, 689, 690], [689, 690, 691], [690, 691, 692], [691, 692, 693], [692, 693, 694], [693, 694, 695], [694, 695, 696], [695, 696, 697], [696, 697, 698], [697, 698, 699], [698, 699, 700], [700, 701, 702], [701, 702, 703], [702, 703, 704], [703, 704, 705], [704, 705, 706], [705, 706, 707], [706, 707, 708], [707, 708, 709], [708, 709, 710], [709, 710, 711], [710, 711, 712], [711, 712, 713], [712, 713, 714], [713, 714, 715], [714, 715, 716], [715, 716, 717], [716, 717, 718], [717, 718, 719], [718, 719, 720], [719, 720, 721], [720, 721, 722], [721, 722, 723], [722, 723, 724], [723, 724, 725], [724, 725, 726], [725, 726, 727], [726, 727, 728], [727, 728, 729], [728, 729, 730], [729, 730, 731], [730, 731, 732], [731, 732, 733], [732, 733, 734], [733, 734, 735], [734, 735, 736], [735, 736, 737], [736, 737, 738], [737, 738, 739], [738, 739, 740], [739, 740, 741], [740, 741, 742], [741, 742, 743], [742, 743, 744], [743, 744, 745], [744, 745, 746], [745, 746, 747], [746, 747, 748], [747, 748, 749], [748, 749, 750], [749, 750, 751], [750, 751, 752], [751, 752, 753], [752, 753, 754], [753, 754, 755], [754, 755, 756], [755, 756, 757], [756, 757, 758], [757, 758, 759], [758, 759, 760], [759, 760, 761], [760, 761, 762], [761, 762, 763], [762, 763, 764], [763, 764, 765], [764, 765, 766], [765, 766, 767], [766, 767, 768], [767, 768, 769], [768, 769, 770], [769, 770, 771], [770, 771, 772], [771, 772, 773], [772, 773, 774], [773, 774, 775], [774, 775, 776], [775, 776, 777], [776, 777, 778], [777, 778, 779], [778, 779, 780], [779, 780, 781], [780, 781, 782], [781, 782, 783], [782, 783, 784], [783, 784, 785], [784, 785, 786], [785, 786, 787], [786, 787, 788], [787, 788, 789], [788, 789, 790], [789, 790, 791], [790, 791, 792], [791, 792, 793], [792, 793, 794], [793, 794, 795], [794, 795, 796], [795, 796, 797], [796, 797, 798], [797, 798, 799], [798, 799, 800], [800, 801, 802], [801, 802, 803], [802, 803, 804], [803, 804, 805], [804, 805, 806], [805, 806, 807], [806, 807, 808], [807, 808, 809], [808, 809, 810], [809, 810, 811], [810, 811, 812], [811, 812, 813], [812, 813, 814], [813, 814, 815], [814, 815, 816], [815, 816, 817], [816, 817, 818], [817, 818, 819], [818, 819, 820], [819, 820, 821], [820, 821, 822], [821, 822, 823], [822, 823, 824], [823, 824, 825], [824, 825, 826], [825, 826, 827], [826, 827, 828], [827, 828, 829], [828, 829, 830], [829, 830, 831], [830, 831, 832], [831, 832, 833], [832, 833, 834], [833, 834, 835], [834, 835, 836], [835, 836, 837], [836, 837, 838], [837, 838, 839], [838, 839, 840], [839, 840, 841], [840, 841, 842], [841, 842, 843], [842, 843, 844], [843, 844, 845], [844, 845, 846], [845, 846, 847], [846, 847, 848], [847, 848, 849], [848, 849, 850], [849, 850, 851], [850, 851, 852], [851, 852, 853], [852, 853, 854], [853, 854, 855], [854, 855, 856], [855, 856, 857], [856, 857, 858], [857, 858, 859], [858, 859, 860], [859, 860, 861], [860, 861, 862], [861, 862, 863], [862, 863, 864], [863, 864, 865], [864, 865, 866], [865, 866, 867], [866, 867, 868], [867, 868, 869], [868, 869, 870], [869, 870, 871], [870, 871, 872], [871, 872, 873], [872, 873, 874], [873, 874, 875], [874, 875, 876], [875, 876, 877], [876, 877, 878], [877, 878, 879], [878, 879, 880], [879, 880, 881], [880, 881, 882], [881, 882, 883], [882, 883, 884], [883, 884, 885], [884, 885, 886], [885, 886, 887], [886, 887, 888], [887, 888, 889], [888, 889, 890], [889, 890, 891], [890, 891, 892], [891, 892, 893], [892, 893, 894], [893, 894, 895], [894, 895, 896], [895, 896, 897], [896, 897, 898], [897, 898, 899], [898, 899, 900], [900, 901, 902], [901, 902, 903], [902, 903, 904], [903, 904, 905], [904, 905, 906], [905, 906, 907], [906, 907, 908], [907, 908, 909], [908, 909, 910], [909, 910, 911], [910, 911, 912], [911, 912, 913], [912, 913, 914], [913, 914, 915], [914, 915, 916], [915, 916, 917], [916, 917, 918], [917, 918, 919], [918, 919, 920], [919, 920, 921], [920, 921, 922], [921, 922, 923], [922, 923, 924], [923, 924, 925], [924, 925, 926], [925, 926, 927], [926, 927, 928], [927, 928, 929], [928, 929, 930], [929, 930, 931], [930, 931, 932], [931, 932, 933], [932, 933, 934], [933, 934, 935], [934, 935, 936], [935, 936, 937], [936, 937, 938], [937, 938, 939], [938, 939, 940], [939, 940, 941], [940, 941, 942], [941, 942, 943], [942, 943, 944], [943, 944, 945], [944, 945, 946], [945, 946, 947], [946, 947, 948], [947, 948, 949], [948, 949, 950], [949, 950, 951], [950, 951, 952], [951, 952, 953], [952, 953, 954], [953, 954, 955], [954, 955, 9
```


RESULTS

Markov Chain

Positive vs Negative samples of Wikipedia with COCA

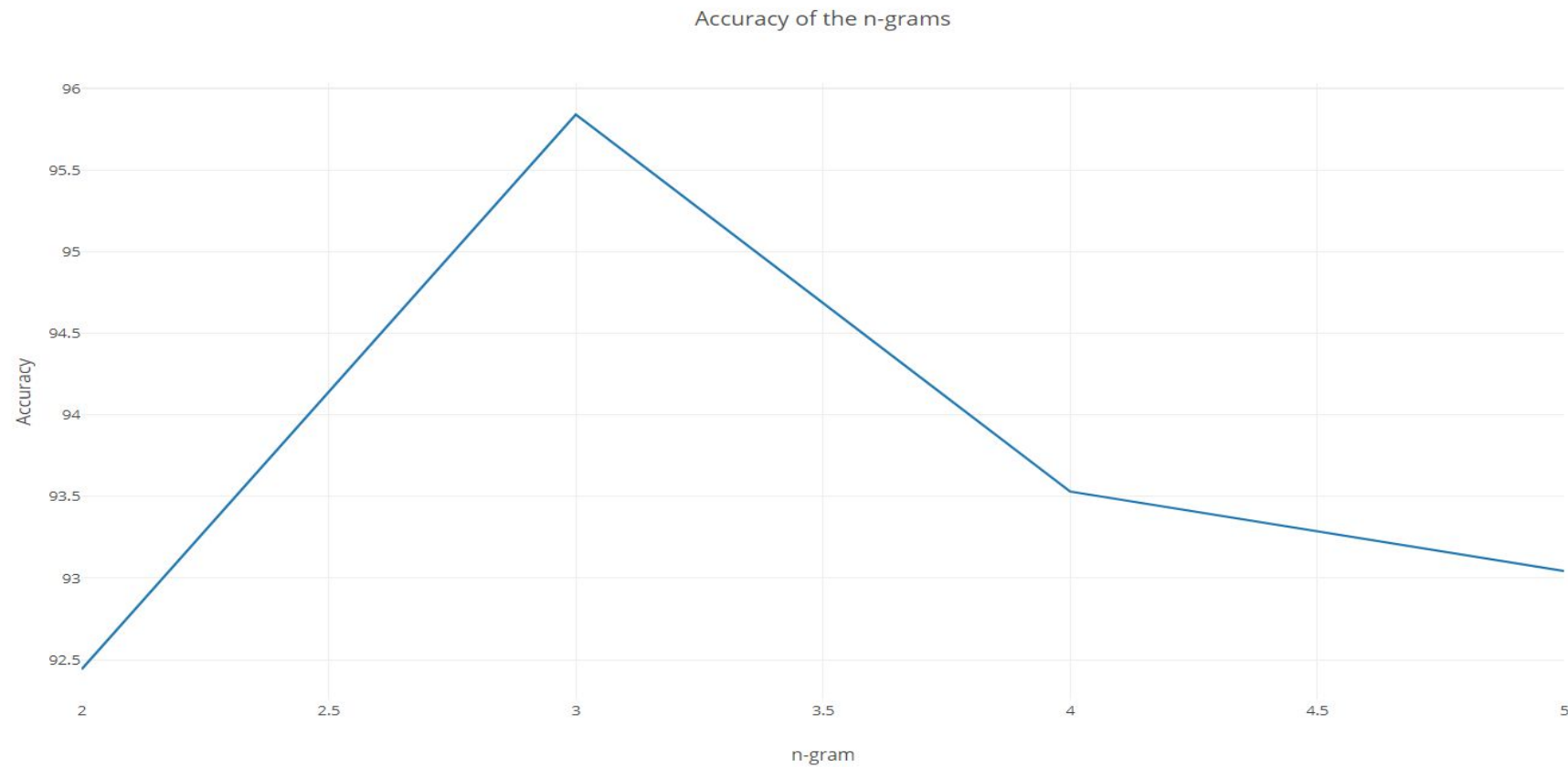


Coca tested over Wikipedia

Table 1 Classification Accuracy and Error Rate for Wikipedia Corpus on COCA Dataset

N-gram	classification accuracy	Error Rate
2	92.44	7.6
3	95.84	4.16
4	93.53	6.47
5	93.042	6.958

n-gram vs accuracy



Spoken and Written American English dataset accuracy

Table 2 Spoken and Written American English Dataset Accuracy

N-gram	w_acad	w_fic	w_mag	w_news	w_spok
2	67.89	62.83	66.74	68.91	67.84
3	72.97	69.27	68.76	71.98	70.29
4	71.67	66.54	67.45	70.27	69.92
5	70.23	67.78	67.23	69.89	69.02

sameerg@sameerg-Inspiron-3543: ~/SEM_6/NLP/Project/Markov

Wi-Fi En Bluetooth (70%) 7:46 PM

```
OrderedDict([('for', 23),
('here', 20),
('this', 19),
('touch', 17),
('such', 16),
('or', 15),
('any', 14)])
```

favor of a new york city-based

sameerg@sameerg-Inspiron-3543:~/SEM_6/NLP/Project/Markov\$ python3 markov.py -w database/big/w4.db -n 4 -words 100 --predict are you in

```
OrderedDict([('the', 114),
('a', 81),
('favor', 67),
('love', 54),
('there', 48),
('trouble', 30),
('pain', 25),
('your', 24),
('for', 23),
('here', 20),
('this', 19),
('touch', 17),
('such', 16),
('or', 15),
('any', 14)])
```

a minute and a

sameerg@sameerg-Inspiron-3543:~/SEM_6/NLP/Project/Markov\$ python3 markov.py -w database/big/w4.db -n 4 -words 100 --predict are you in

```
OrderedDict([('the', 114),
('a', 81),
('favor', 67),
('love', 54),
('there', 48),
('trouble', 30),
('pain', 25),
('your', 24),
('for', 23),
('here', 20),
('this', 19),
('touch', 17),
('such', 16),
('or', 15),
('any', 14)])
```

the face and neck

sameerg@sameerg-Inspiron-3543:~/SEM_6/NLP/Project/Markov\$

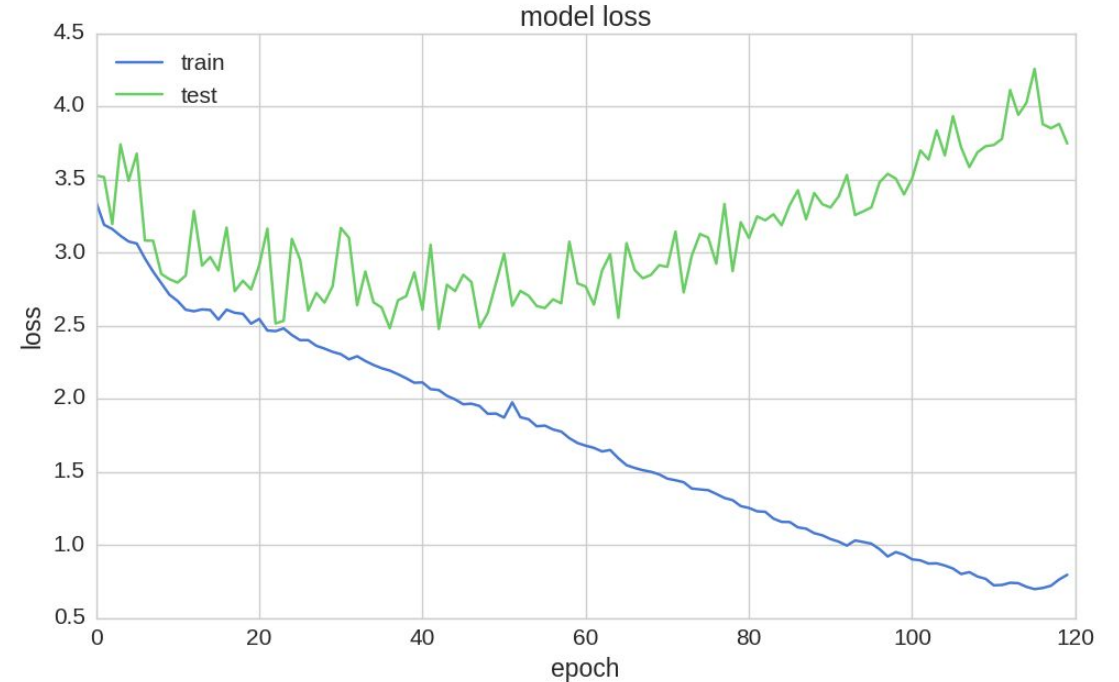
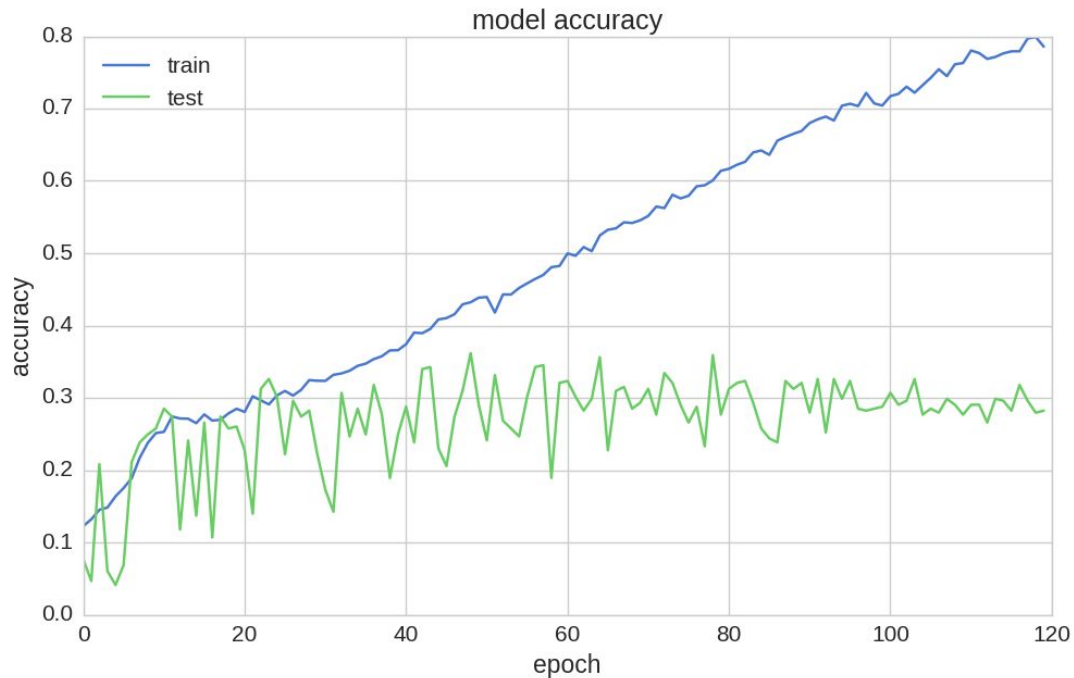
Outputs:

```
g@sameerg-Inspiron-3543: ~/SEM_6/NLP/Project/Markov
('committed', 26),
('denying', 26),
('friends', 26),
('talkin', 26),
('already', 25),
('any', 25),
('proposing', 25),
('tonight', 25)])
saying you have the same thing with me and my father had a chance for the last time i saw a man in the middle east peace conference on the floor of the most part they
sameerg@sameerg-Inspiron-3543:~/SEM_6/NLP/Project/Markov$ python3 markov.py -w database/big/w3.db -n 3 -words 100 --predict are you
OrderedDict([('going', 7993),
('doing', 6047),
('saying', 2096),
('talking', 2036),
('sure', 1942),
('a', 1524),
('all', 1235),
('ready', 1150),
('in', 924),
('still', 858),
('kidding', 824),
('okay', 729),
('looking', 709),
('trying', 699),
('there', 654),
('feeling', 613),
('getting', 562),
('committed', 26),
('denying', 26),
('friends', 26),
('talkin', 26),
('already', 25),
('any', 25),
('proposing', 25),
('tonight', 25)])
saying the words that are in a moment of silence in which the united states is the most important things to do it again and again to the point of view and the united nations and other things to say that the government to the united kingdom and the united kingdom and the other day that i have n't had the same way i was a very different than what you want me to do it for you in your life in the world and to the point is that the president and ceo john
sameerg@sameerg-Inspiron-3543:~/SEM_6/NLP/Project/Markov$
```

output for 100 words given a bigram , and trigram .database w3.db is used. with starting by predicting the 3 gram , then last two words of the tri gram is used to predict the next 3-gram till 100 words are produced.

RNN using character approach:

First we trained by predicting character by character and achieved an accuracy of 78.57% over 120 epochs but then, this didn't suit our application and results weren't that satisfying



As we can see , here the results aren't that pleasing when seen with the actual sentence in the quotes list

```
In [42]: quotes = ["Mr. Dursley was the director of a firm called Grunnings, which made drills. He was a big, beefy man with hardly  
"At half past eight, Mr. Dursley picked up his briefcase, pecked Mrs. Dursley on the cheek, and tried to kiss Dudley good-  
,"he'd gotten into terrible trouble for being found on the roof of the school kitchens"]
```

```
In [43]: for q in quotes:  
    seq = q[:40].lower()  
    print(seq)  
    print(predict_completions(seq, 5))  
    print()
```

```
mr. dursley was the director of a firm c  
['lout... ', 'er ', 'as ', 'houlded ', 'it ']
```

```
at half past eight, mr. dursley picked u  
['itew ', 'phow ', 'she ', 'moffed ', 'whordind ']
```

```
he'd gotten into terrible trouble for be  
['chouddare. ', 'lhed ', 'theded ', 'eked ', 'kne ']
```

```
In [ ]:
```


Results of LSTM for word by word:

Approach / Dataset	US Speech	GST Speech
3 Word approach	<u>87.28%</u>	<u>80.33%</u>
Whole sentence approach	<u>96.28%</u>	<u>95.60%</u>

Model Accuracy on Modi's Speeches after 150 epochs

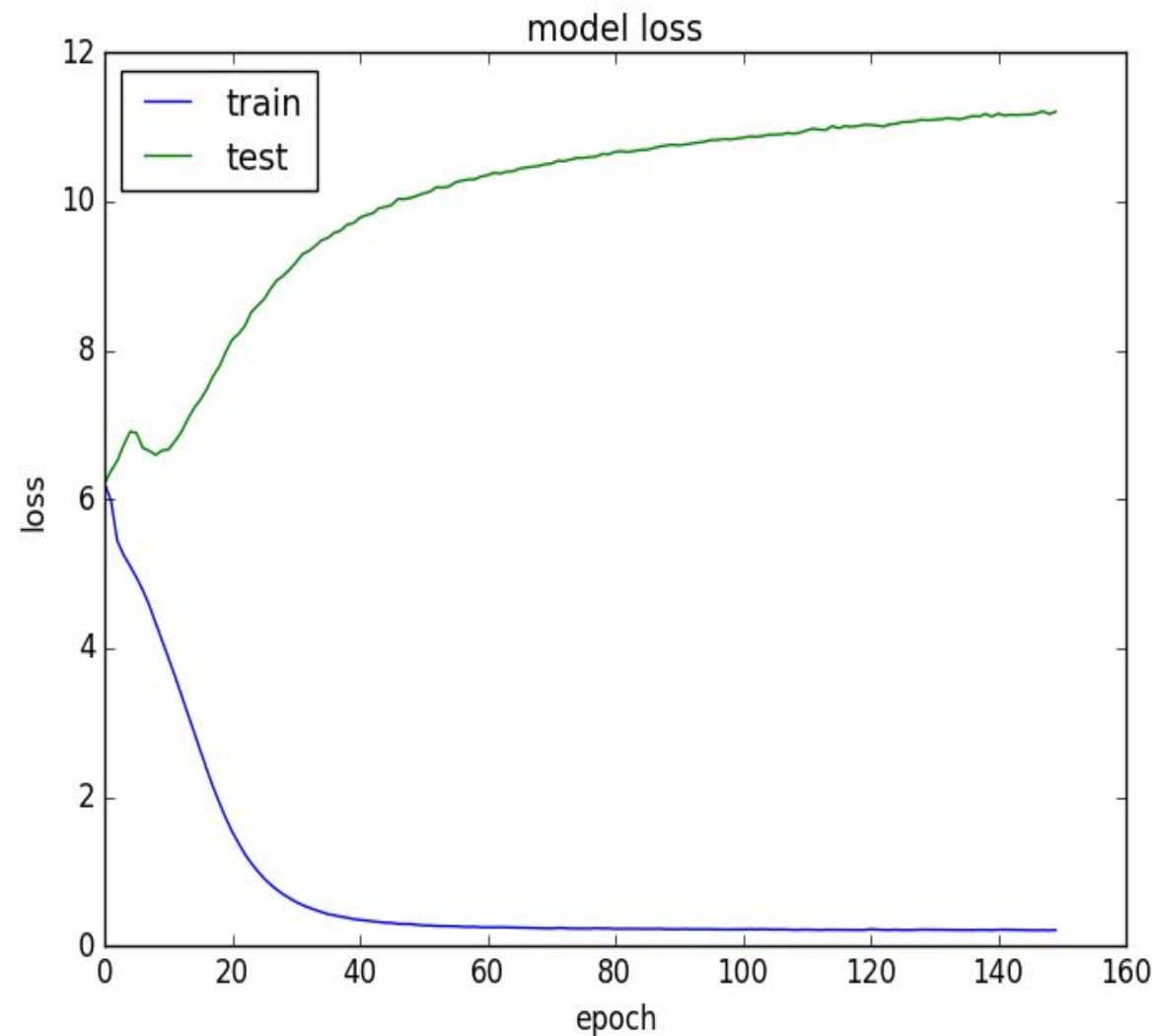
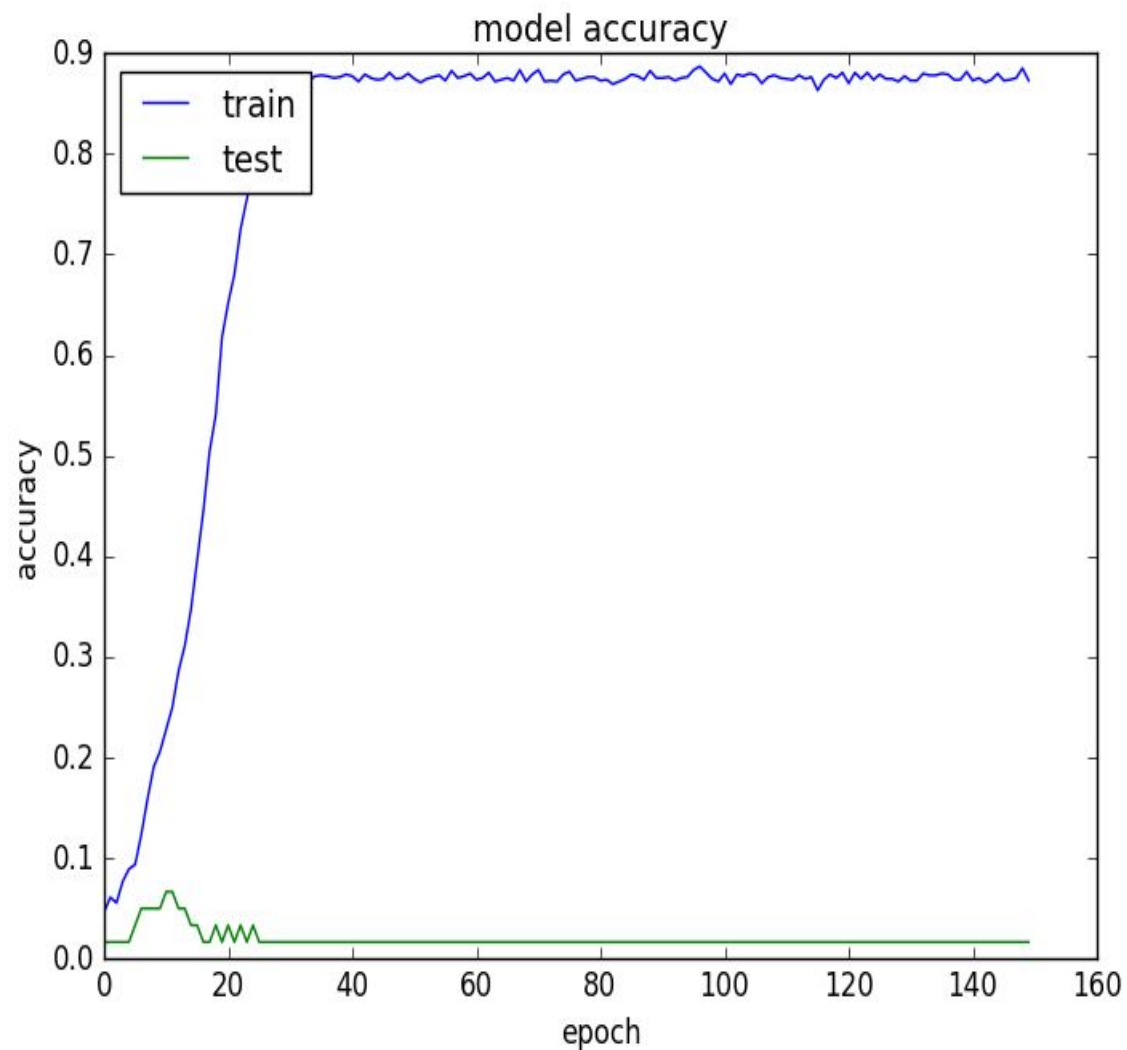
Model Accuracies (Adam VS RMSprop) for Harry Potter Novel after 120 epochs

Approach / Optimizer	Adam	RMSprop
3 Word approach	<u>56.44%</u>	<u>37.74%</u>

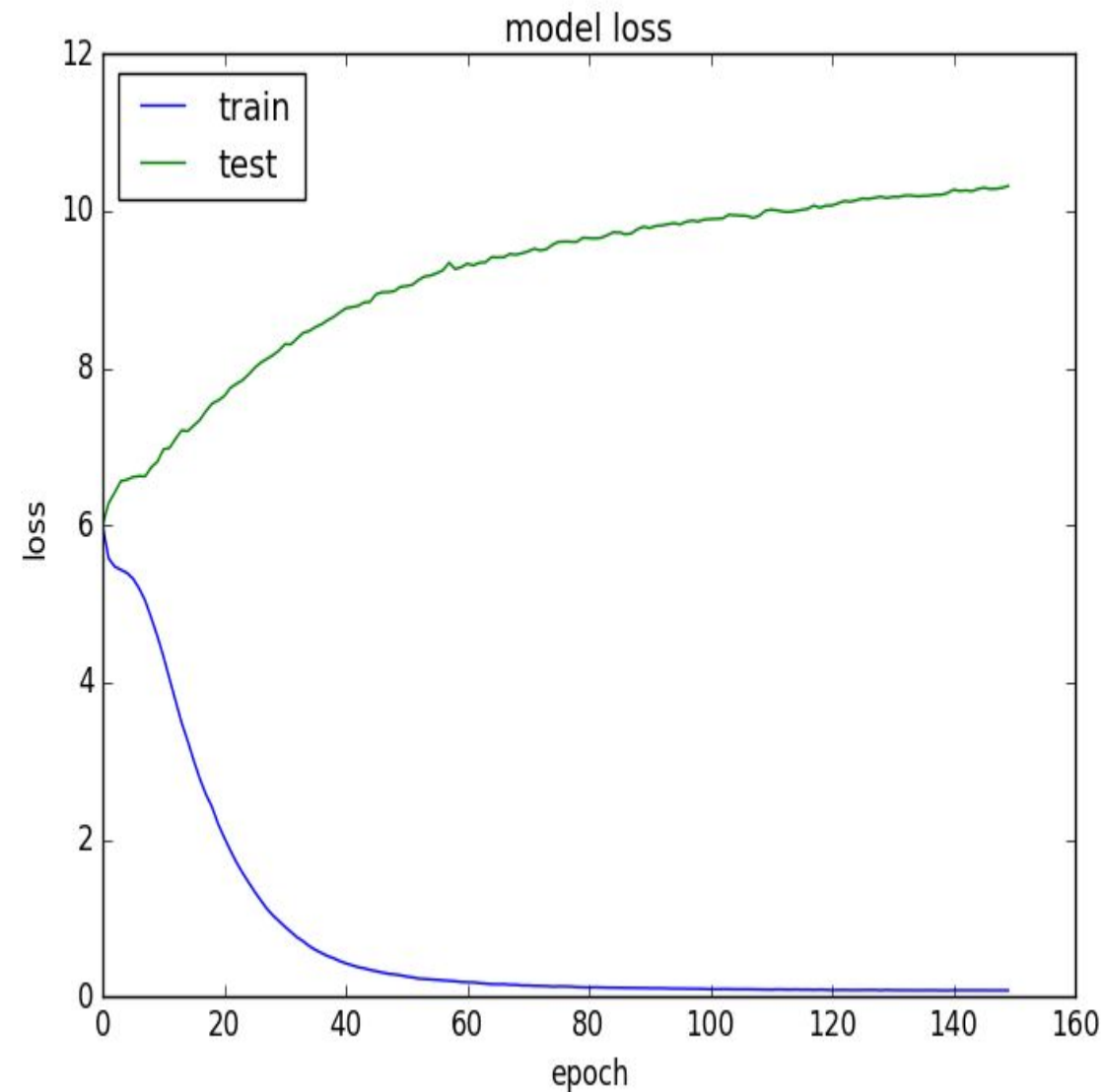
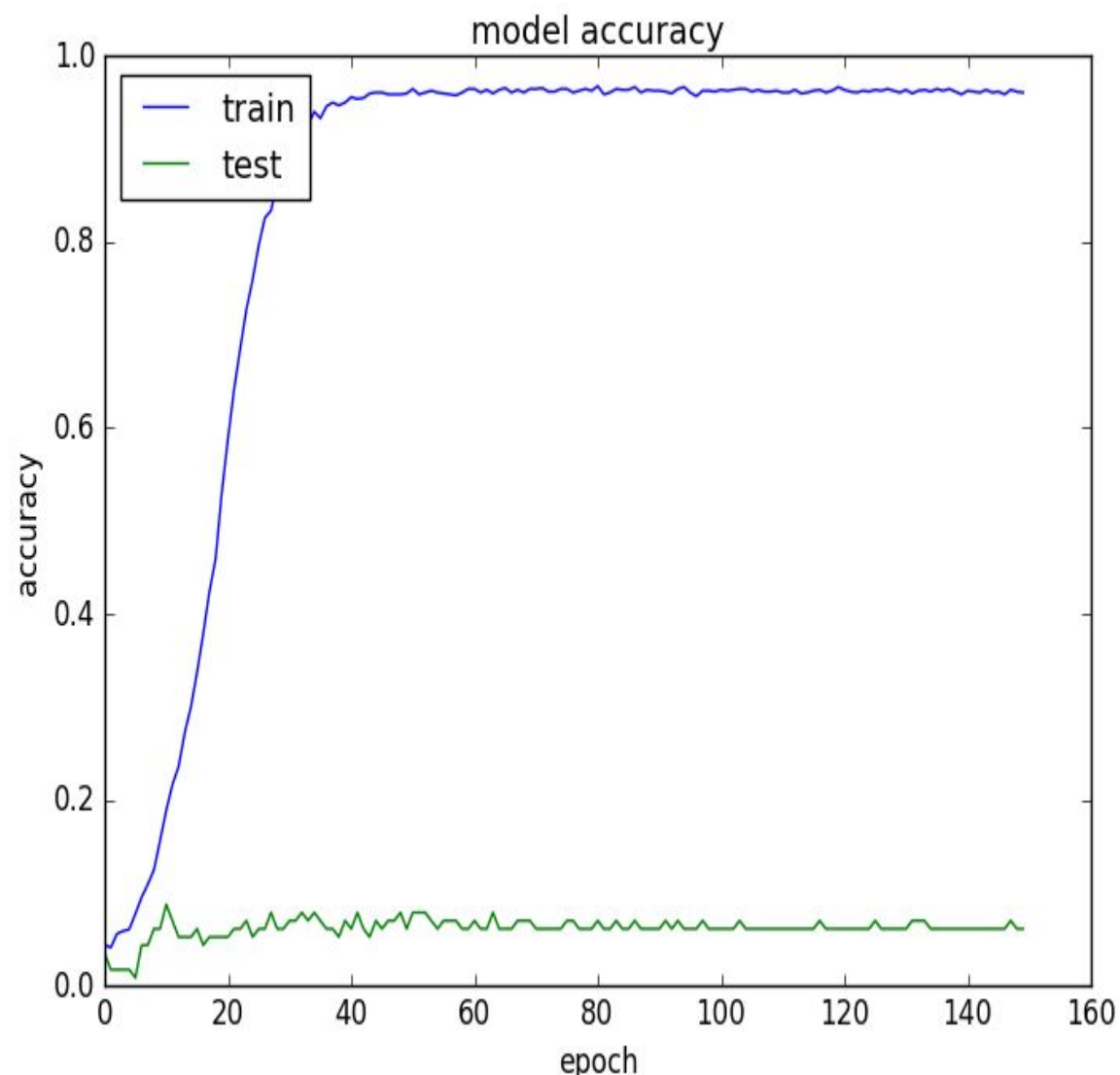


Epochs vs Accuracy/loss :

3 Word - US Speech



Sentence - US Speech




```
print(generate_seq(loader_model, tokenizer, 2, 'people', 2))
```

people in the

```
print(generate_seq(loader_model, tokenizer, 2, 'United Nations', 50))
```

United Nations including its security council so that it carries greater credibility and legitimacy and will be more representative and effective in achieving our goals are comprehensive it gives priority to the future of the island states i speak about blue revolution which includes the prosperity sustainable use of marine wealth and

Fig: Predictions on 3 Word Model for US Speech

```
print(generate_seq(loader_model, tokenizer, 57, 'People', 2))
```

People in the

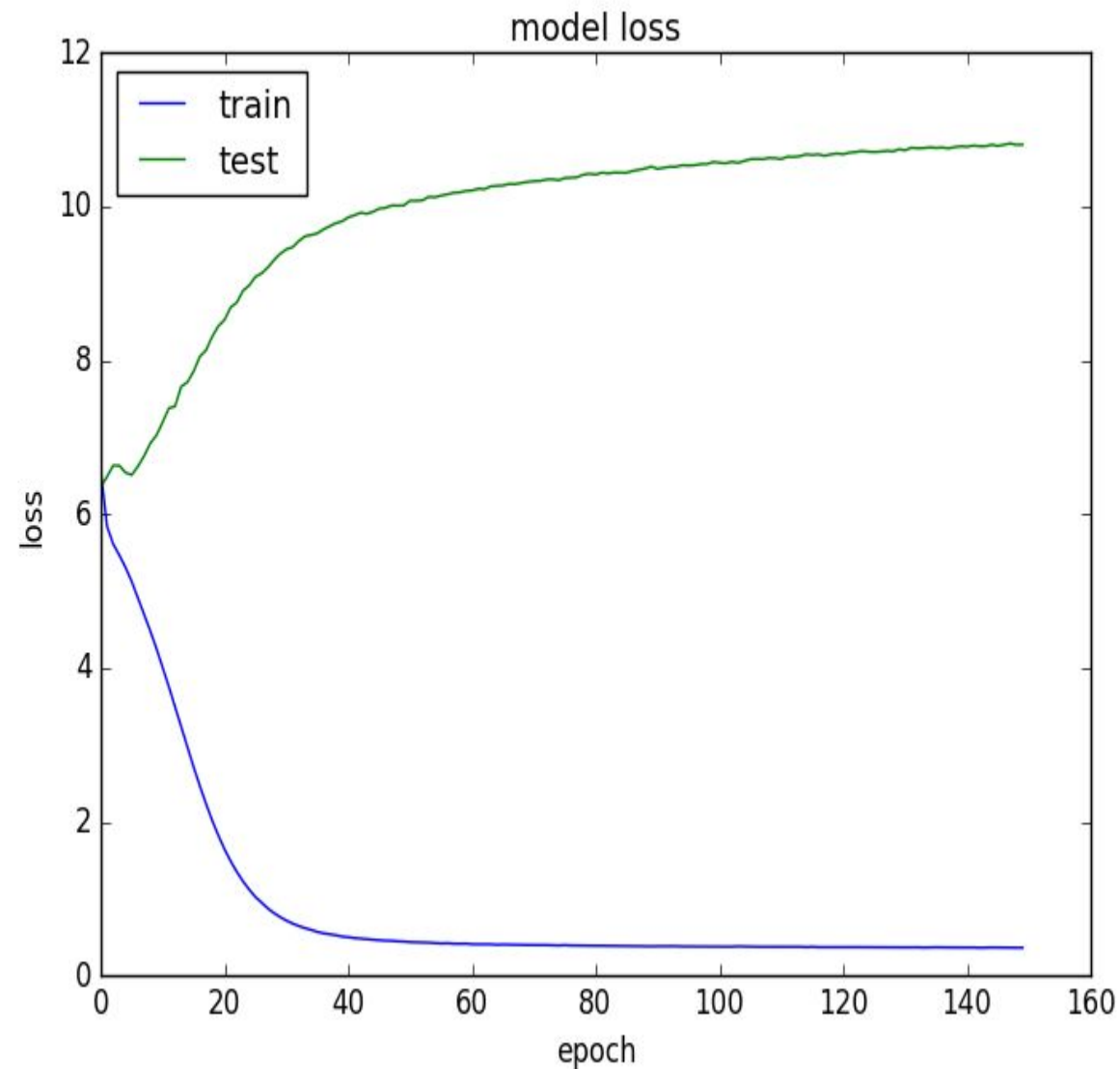
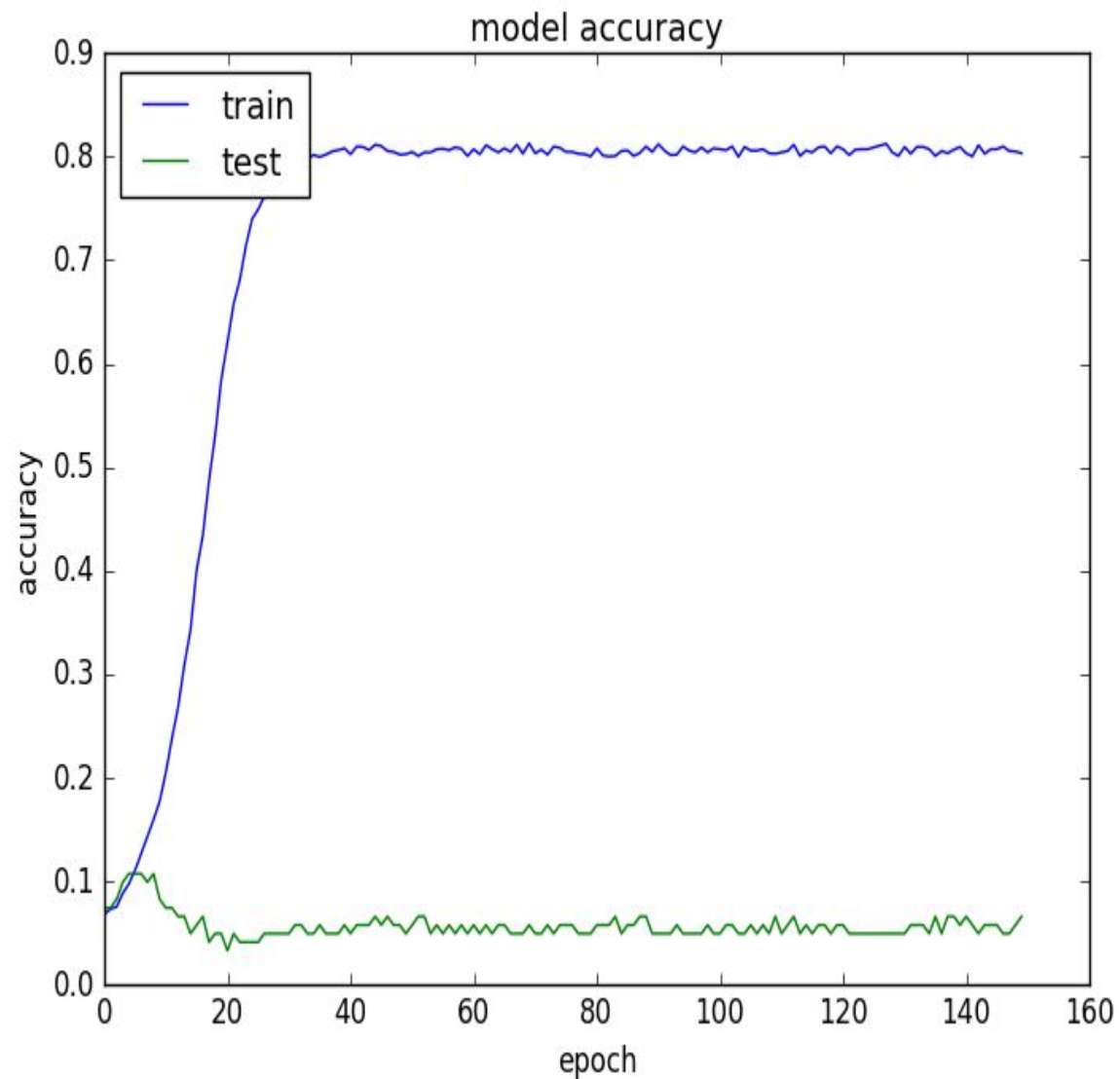
```
print(generate_seq(loader_model, tokenizer, 57, 'United Nations', 50))
```

United Nations see now distance is no insulation from challenges and privations from distant lands states from the pacific to the atlantic of natural disasters and pension for everyone's sunset years energy efficiency a tax on coal a huge afforestation programme reforming our transportation and cleaning up our cities and rivers up

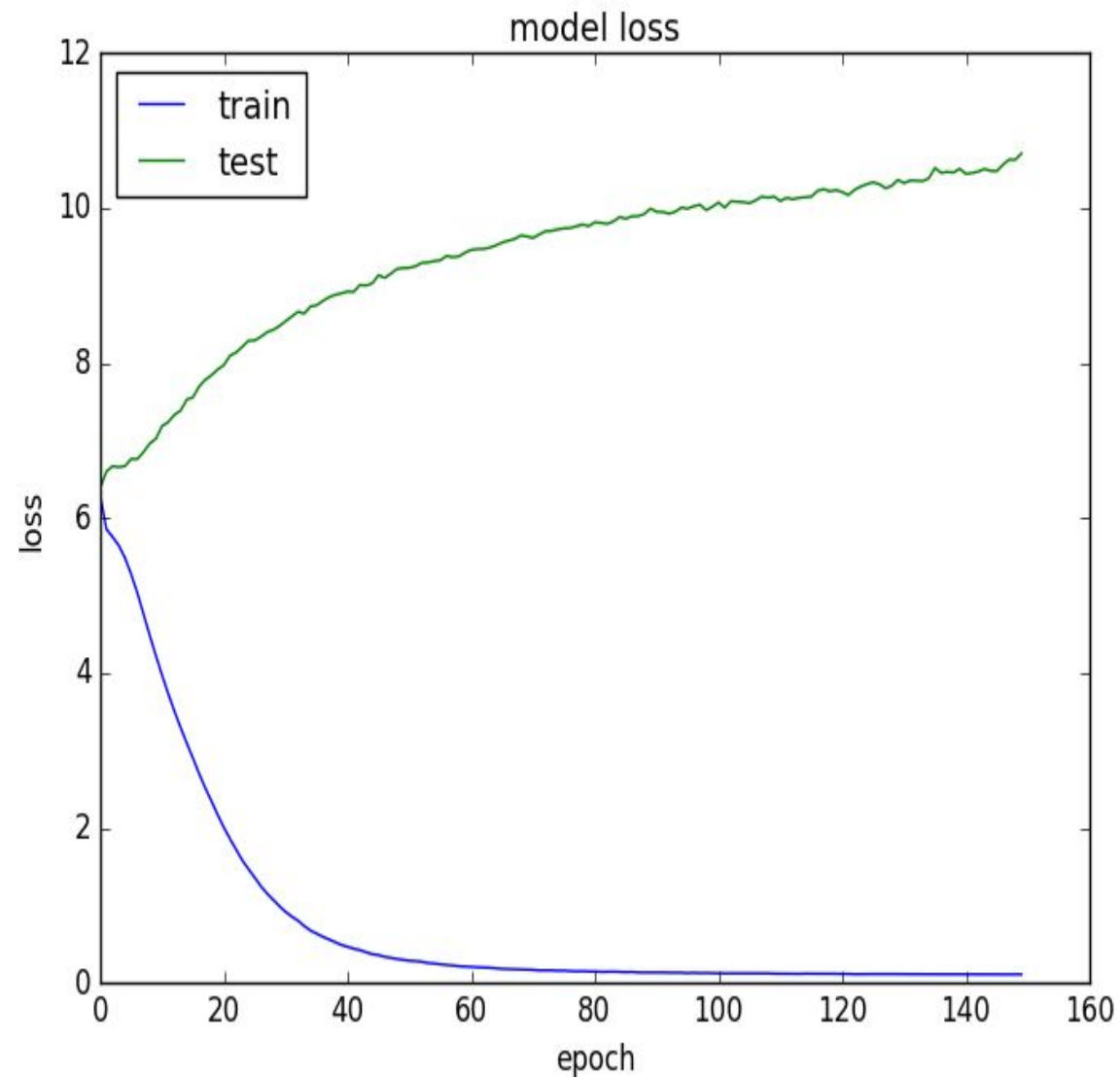
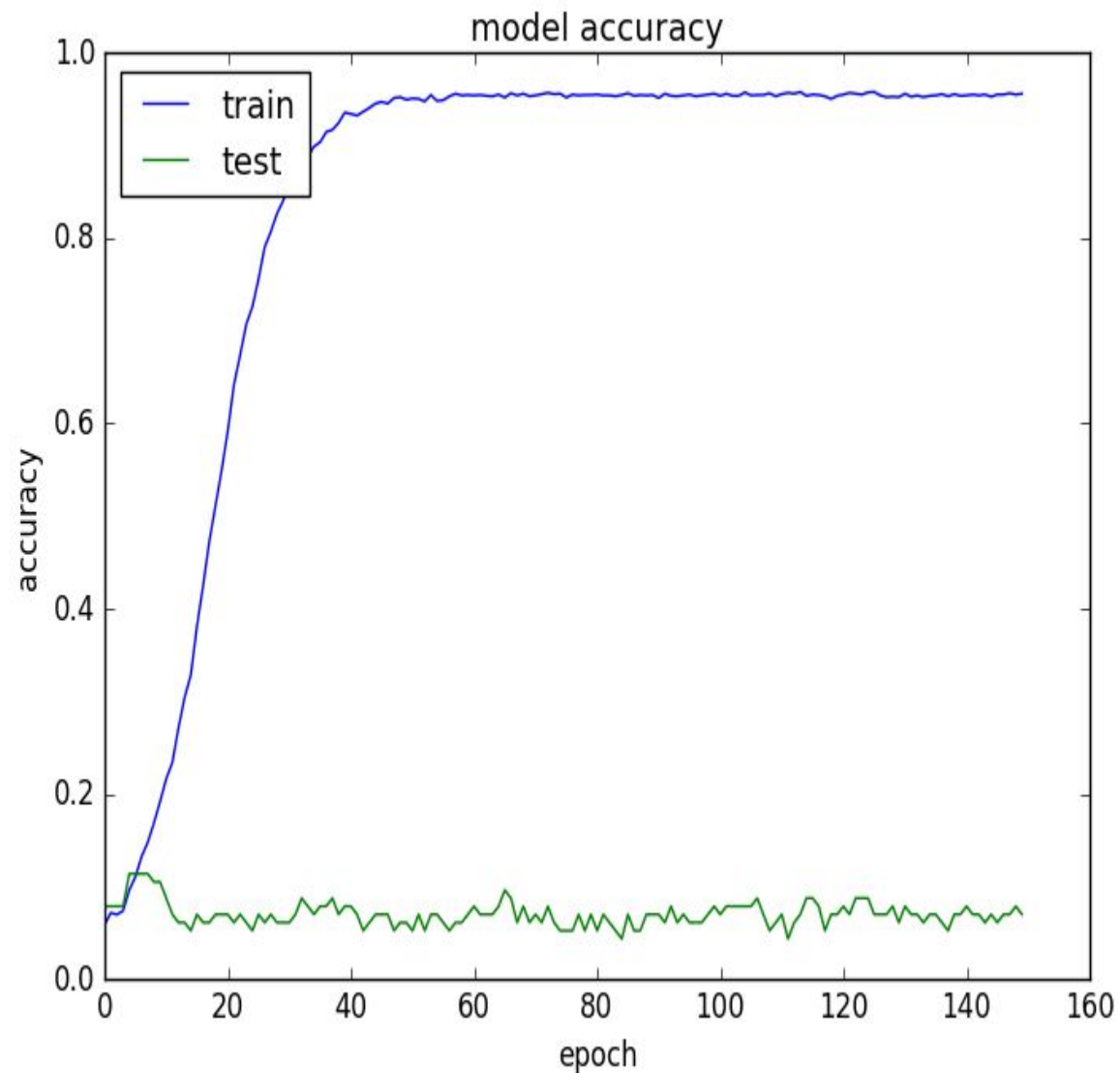
Fig: Predictions on Sentence Model for US Speech

We see that both the models have pretty much given a meaningful set of words as output. But the sentence approach is close the actual data but the three word trade off looks very meaningful too.

3 Word - GST



Sentence - GST




```
print(generate_seq(loader_model, tokenizer, 2, 'important', 50))
```

important when the constitution was framed it introduced a framework of equal opportunity for development gst is the dream of our country's poor as we expand horizontally we need to ensure vertical growth as well the work in unison for india under gst for the first time the central service officers

```
print(generate_seq(loader_model, tokenizer, 2, 'State Governments', 100))
```

State Governments there have been done for this new system but those arrangements are at basic level and the esteemed conglomerate gathered here representing diverse sectors in the same product will have different prices in say delhi gurugram and noida which are 25 30 kilometres apart this is a landmark step towards economic reforms beyond the taxation revamp it is a landmark step towards economic reforms beyond the taxation revamp it is a landmark step towards economic reforms beyond the taxation revamp it is a landmark step towards economic reforms

Fig: Predictions on 3 Word Model for Speech on GST

```
print(generate_seq(loader_model, tokenizer, 66, 'important', 50))
```

important when the constitution was framed it introduced a framework of equal opportunity and rights for all those troubles as the poor so that the poor is benefitted just in a india of one of the biggest strength of federal structure the gst reform be on the nation and accepted by

```
print(generate_seq(loader_model, tokenizer, 66, 'State Governments', 100))
```

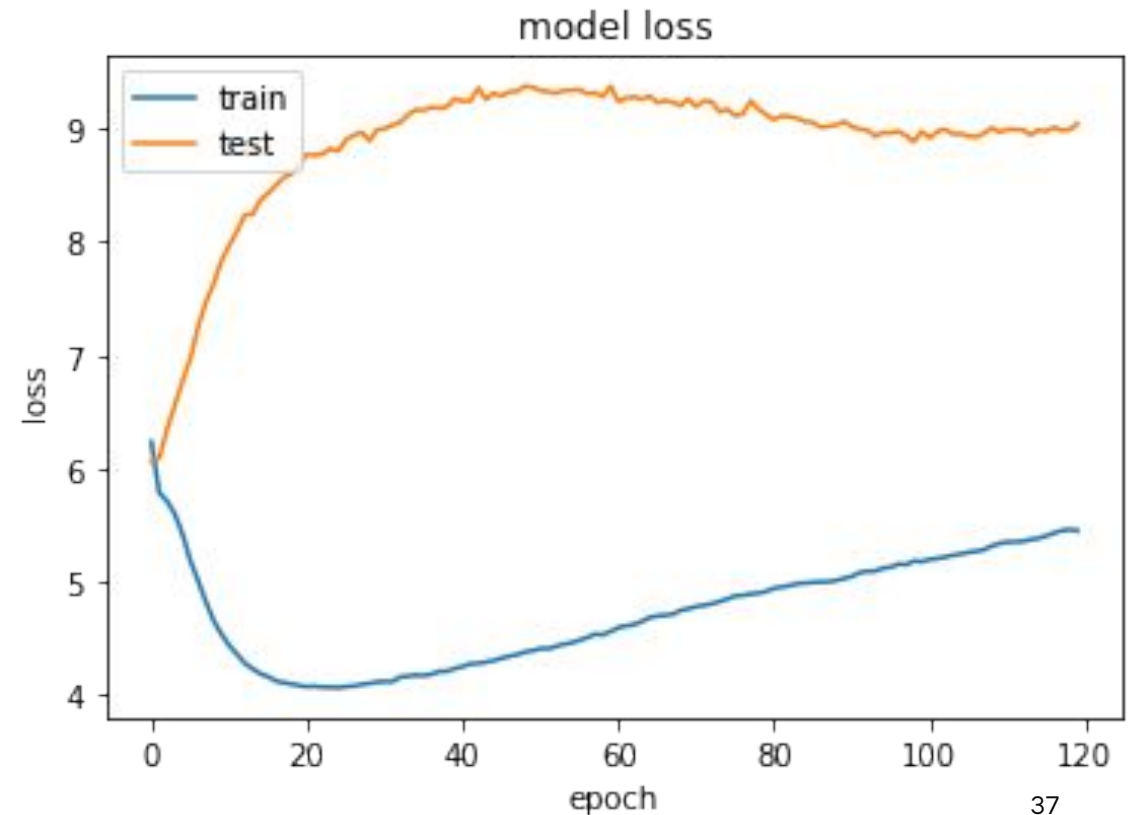
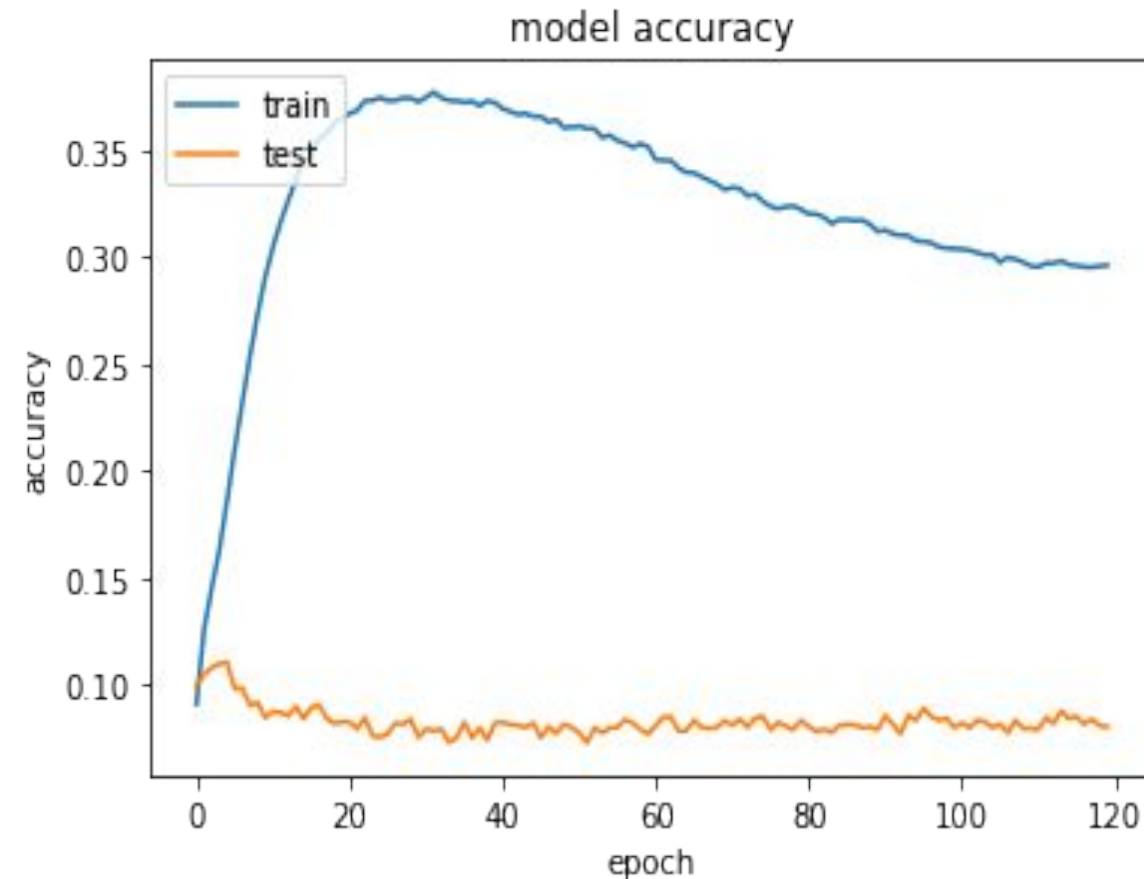
State Governments will be one nation one tax and tax system shall be executed in a standard manner in all the states and they are slated to gain immensely as they shall now get equal opportunities of development and the constitution was accepted by the nation and they will be shred off to have one nation one tax right from ganga nagar to itanagar and from leh to enriched take take it resulted in take the example of the country participated in the country's constitution was accepted by the nation attaining freedom on 26th november 1949 this house stands as a testimony to

Fig: Predictions on Sentence Model for Speech on GST

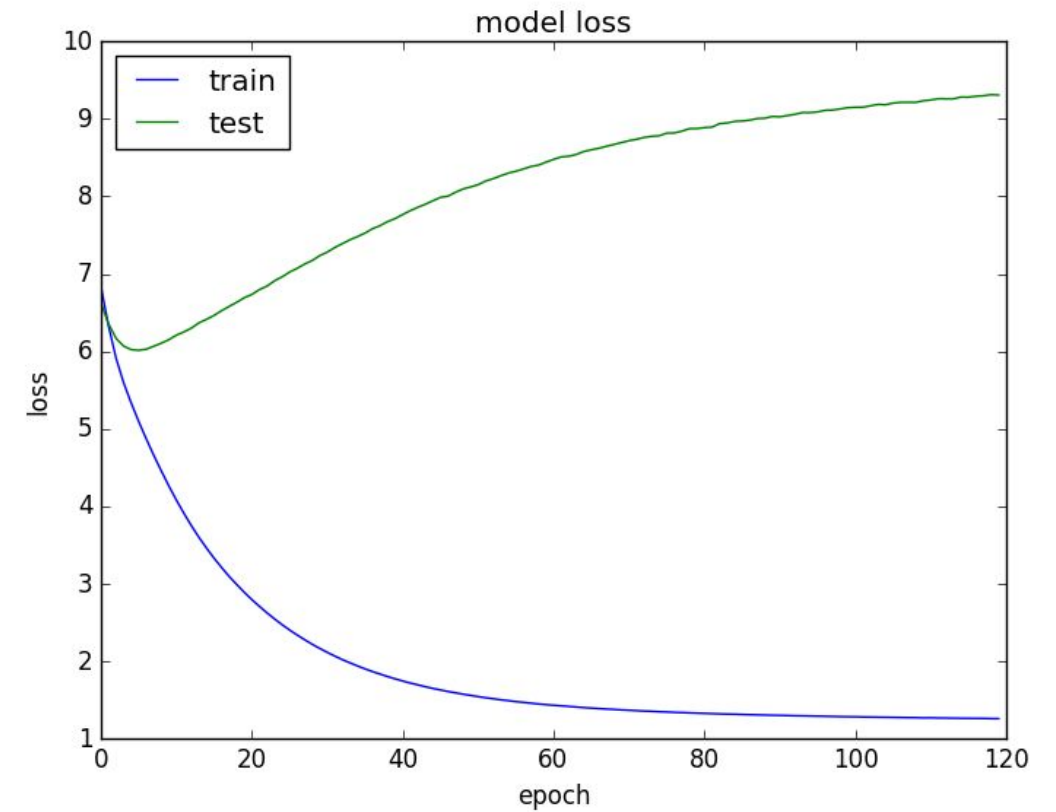
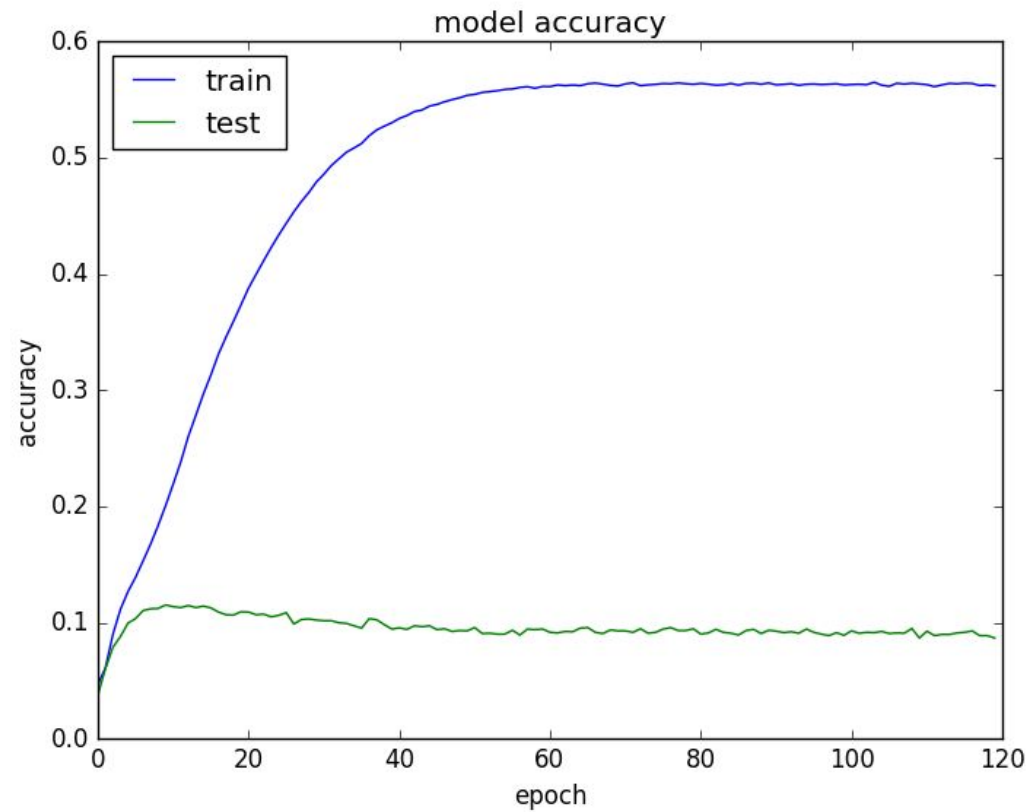
Even here , we see that both the models have pretty much given a meaningful set of words as output. But the sentence approach is close the actual data but the three word trade off looks very meaningful too. ³⁶

RMSprop V/s Adam for Harry Potter Novel

RMSprop optimization fetched an best accuracy of 37.74%



Adam fetched an accuracy of 56.44%, which even has a better validation accuracy when compared to RMSprop




```
print(generate_seq(loader_model, tokenizer, 2, 'hogwarts woke', 50))
```

hogwarts woke to go and once and harry was there he was going to be a back to the dursleys had been a back to the dursley
s had been a back to the dursleys had been a back to the dursleys had been a back to the dursleys had been a back

```
print(generate_seq(loader_model, tokenizer, 2, 'It took quite a while for them all to get off', 10))
```

It took quite a while for them all to get off the morning it was a very good term said dumbledore

FIG: Test Outputs on Harry Potter Novel using RMSprop optimizer

```
print(generate_seq(model, tokenizer, max_length-1, 'hogwarts woke', 50))
```

hogwarts woke to find out how to get past fluffy aren't you wearing yours ron george demanded come on harry an
d ron thought see professor snape disliked him by his proper name voldemort professor mcgonagall was out of th
e way to the door and said she had a long time he was

```
print(generate_seq(model, tokenizer, max_length-1, 'It took quite a while for them all to get off', 10))
```

It took quite a while for them all to get off neville for neville had been a mistake said uncle vernon

FIG: Test Outputs on Harry Potter Novel using Adam optimizer

We can see clearly that adam optimization performed well than RMSprop which has repeated words, than adam which actually produced a paragraph with enough sense

Markov Chains V/s RNN

- The RNN is more accurate than the n-gram model because it can retain more history than just the previous 2 tokens and using word vectors instead of discrete tokens for each word allows it to generalize and make predictions from similar words, not just exact matches.
- The RNN takes longer to train (e.g. 8 hours vs 8 seconds, about 4000x longer), but would easily be able to fit into memory for applications like phone text entry - it also performs word predictions a bit slower than the n-gram, but not enough to be noticeable to an end-user.

Character by Character V/s Word by Word LSTM

1. As we can see that the outputs of word by word model were very very better than character by character.
2. And the training time is also less for the word by word model.
3. The two tradeoffs that is the sentence and the three word approaches in word by word models were compared in the previous slides accordingly to the dataset.

Training Times

Dataset/Approach	3 word	Full sentence
US Speech	2 secs/epoch (GEFORCE 920M)	6 secs/epoch (GEFORCE 920M)
GST Speech	~3 secs/epoch(GEFORCE 920M)	18 secs/epoch(GEFORCE 920 M)

Dataset/Approach	3 word/Adam	3 word/RmsProp
HarryPotter	250 secs/epoch (GEFORCE 860M)	194 secs/epoch (GEFORCE 860M)

| MILESTONES

Serial no	Milestone description	Status (% complete)	Comments
1	Markov Chain	100%	Was easy to implement
2	Integrating SQL databases to Markov chain	100%	Was a little tough to maintain the .db files but querying was very fast
3	Writing a simple LSTM network	100%	Initial results were bad
4	Fine tuning the network	80%	We could get accuracy above 90% for NM speech but validation accuracy is less around 15%. If no validation set is given then the results are very good .
5	Have to train the network over large text like novels	100%	We got around 60-70 % accuracy which has to be improved
6	Fine tuning for large text	80%	Still figuring out on the optimizers other than adam to improve the accuracy
7	Dataset for Music compositions	-	Further Work
8	Creating the application for keyboard and music generation -	-	Further work

TOP CHALLENGES UNRESOLVED SO FAR

1. Achieving accuracy above 75% over the large text corpus
2. Improving validation accuracy above 30 % for small text corpus
3. Training time isn't reduced less than 260 sec per epoch on different batch sizes for large text corpus

OUR GOING FORWARD PLAN (OPTIONAL)

1. Step 1 : Creating better architectures to attain better validation accuracy
2. Step 2 : Making an keyboard application to android phones and Survey with people on how good this application is performing with their phones.
3. Step 3: Training the architecture on musical compositions of a particular singer and try to generate predictive musical notes which can uniquely represent his compositions.

OUR TOP THREE LEARNING IN THIS PROJECT

1. Learning # 1 : LSTM property to remember sequences and and produce the next output . This helped us a lot in creating better architectures.
2. Learning # 2 : Markov Chains to predict sequences with their unique property.
3. Learning # 3 : Working with different libraries like nltk , keras ,tensorflow, matplotlib