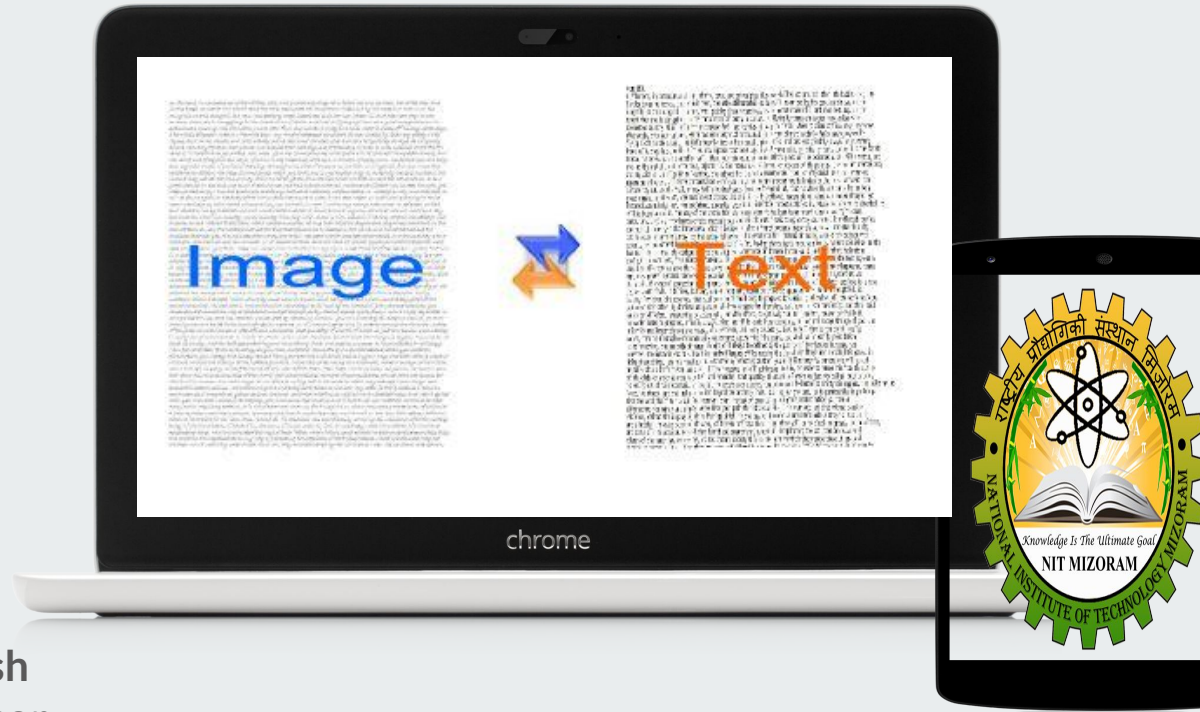# Automatic Image Captioning System

**Saurav Saha**
BT14CS002
7th Semester
Dept. of CSE
NIT MIZORAM

**Supervisor**
**Mr. Sandeep Dash**
Assistant Professor
Dept. of CSE
NIT MIZORAM

# Overview

# INTRODUCTION

- Given an image as input, the system generates automated caption for it.

- Explores the idea of generating embeddings for multiple modalities and project them into same representation space.[1]

- It also helps in either way retrieval as the generated embeddings correspond to each other.

# MOTIVATION

- SYSTEM CAN BE AN AID TO VISUALLY IMPAIRED PEOPLE

- AUTOMATED FRAME-BY-FRAME DESCRIPTION/SUBTITLE GENERATION CAN BE DONE OF VIDEOS

- RELATED IMAGES CAN BE OBTAINED FOR A GIVEN CAPTION, WHICH CAN CONTAIN MORE SEMANTIC INFORMATION

# LITERATURE REVIEW

Prevalent Approaches:

- Template Based
    - Detect objects & attributes
    - Sentence -> Phrase
    - Learn models like CRF
- Retrieval Based
    - Leverage Distance in Visual Space, find image related to test image
    - Combine the caption & modify it
- Neural Networks based[1, 2]
    - Learn Common Embedding
    - Use CNN, RNN, LSTM

## DATA DESCRIPTION

**Insights on DataSet:**

Training Images: 9000

Test Images: 1000

1 caption per image

MS COCO 2014 Validation Dataset[3]

# DATASET



Fig: a baseball player at bat in a game track

Fig: a ski lift carrying people over a snow covered mountain

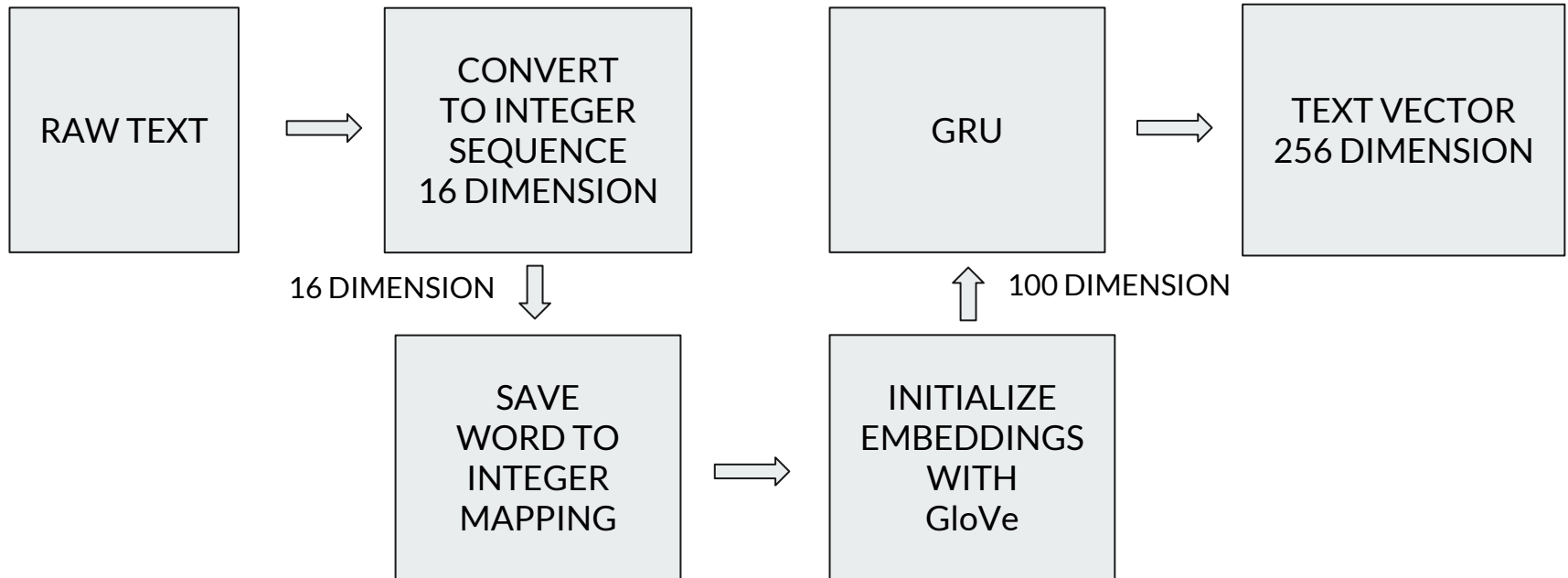Fig: a yellow and blue train on railway track

# SYSTEM DESCRIPTION

- TEXT REPRESENTATION
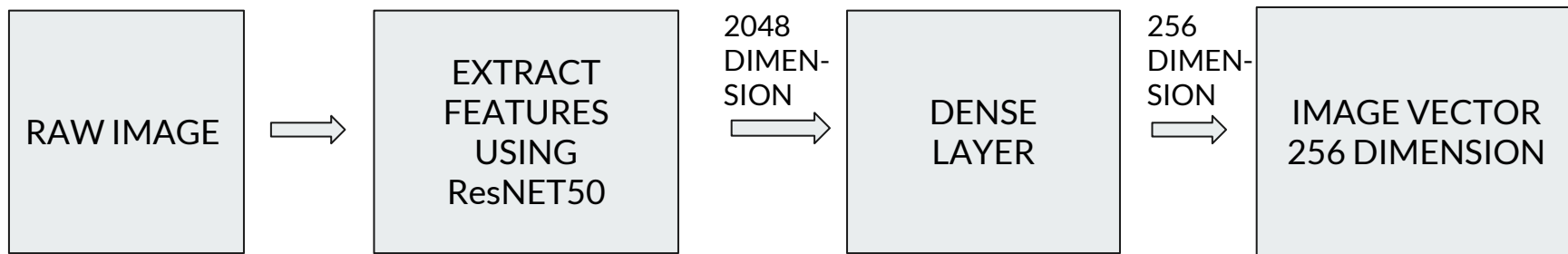- IMAGE REPRESENTATION
- SYSTEM TRAINING
- IMAGE TO CAPTION GENERATION

# TEXT REPRESENTATION

```
RAW TEXT  ⇒  CONVERT
              TO INTEGER
              SEQUENCE
              16 DIMENSION

                16 DIMENSION ⇓

              SAVE          ⇒  INITIALIZE
              WORD TO           EMBEDDINGS
              INTEGER           WITH
              MAPPING           GloVe

                              100 DIMENSION ⇑

              GRU           ⇒  TEXT VECTOR
                              256 DIMENSION
```

# IMAGE REPRESENTATION

RAW IMAGE → EXTRACT FEATURES USING ResNET50 → 2048 DIMENSION → DENSE LAYER → 256 DIMENSION → IMAGE VECTOR 256 DIMENSION

# SYSTEM TRAINING

$$\sum max(0, 1 - pos + neg)$$

**TRAINING PIPELINE**

ACTUAL CAPTION

TEXT VECTOR FOR ACTUAL CAPTION

IMAGE INPUT

IMAGE VECTOR FOR INPUT IMAGE

COMPUTE DOT PRODUCT

POSITIVE PAIR

NEGATIVE PAIR

CONCATENATE

OUTPUT

CUSTOM LOSS

COMPILE TRAINING MODEL

SAVE TRAINED MODEL

NOISE CAPTION

TEXT VECTOR FOR NOISE CAPTION

CUSTOM ACCURACY

$$\mu \ (pos > neg)$$

# GENERATING CAPTION

RAW IMAGE → EXTRACT FEATURES USING ResNET50 → **2048 DIMEN-SION** → DENSE LAYER → **256 DIMEN-SION** → IMAGE VECTOR 256 DIMENSION → COMPUTE DOT PRODUCT → FIND CAPTION → OUTPUT

TEXT VECTORS OF ALL CAPTIONS OF TRAINING SET SAVED AFTER TRAINING MODEL → COMPUTE DOT PRODUCT

# RESULTS

- ANALYSIS ON TRAINING
- SUCCESSFUL CAPTIONS
- PARTIALLY SUCCESSFUL CAPTIONS
- UNSUCCESSFUL CAPTIONS
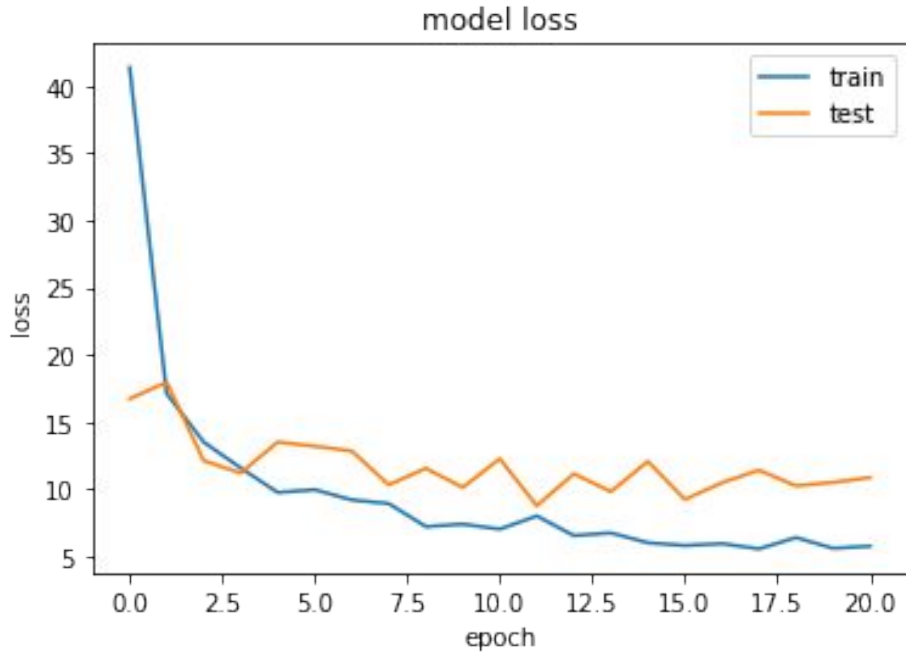- BLEU SCORE

# ANALYSIS ON LOSS


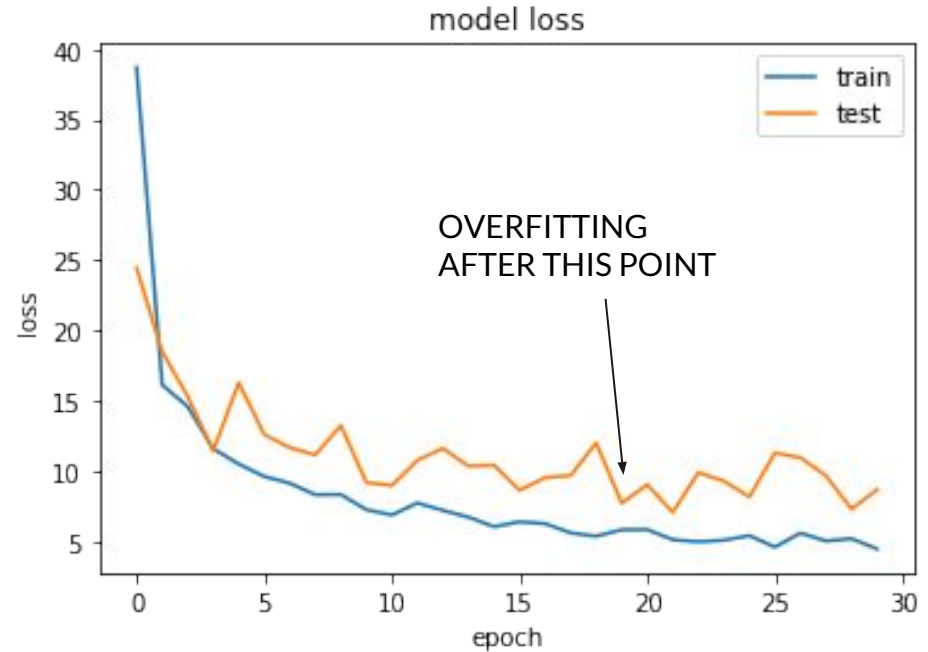
Fig: Loss on 21 epoch

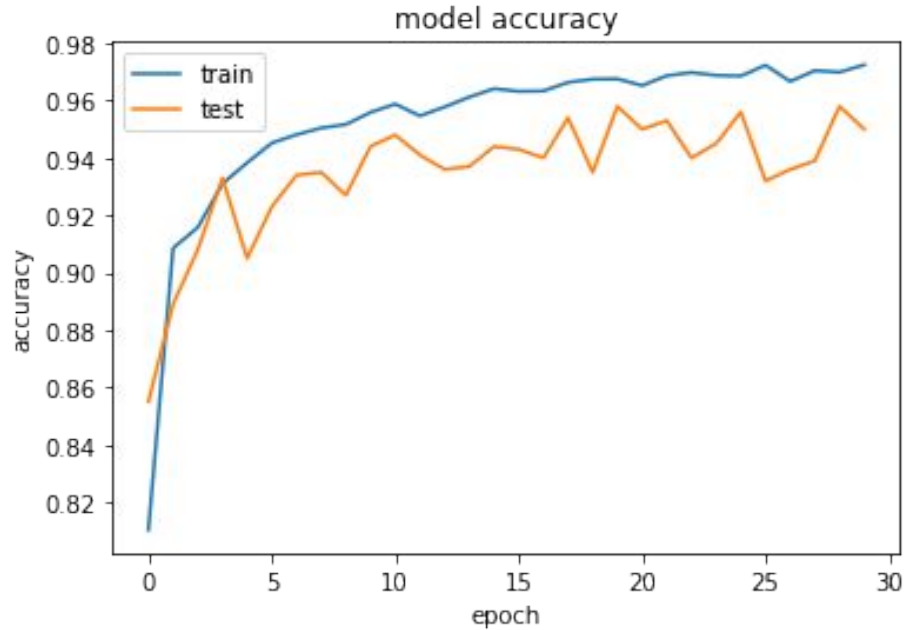Fig: Loss on 30 epoch [OVERFITTING]

# ANALYSIS ON ACCURACY
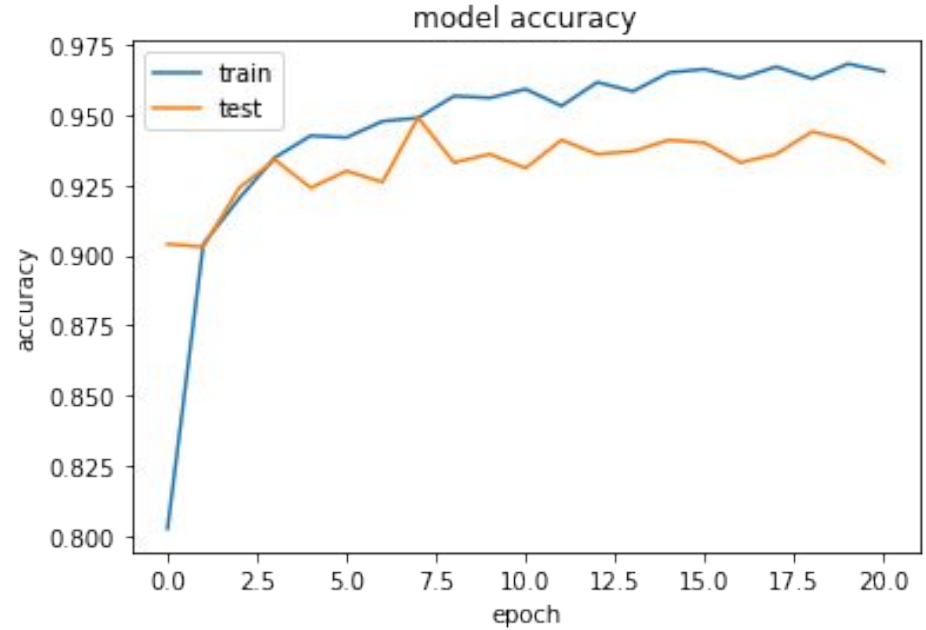


Fig: Accuracy on 30 epoch



Fig: Accuracy on 21 epoch

# SUCCESSFUL CAPTIONS





**Original Caption**: a large long train on a steel track

**System Generated Caption**: a train is stopped at a train station platform

**Original Caption**: this man is skiing down a mountain slope

**System Generated Caption**: a man in skies is walking in the snow

# PARTIALLY SUCCESSFUL CAPTIONS



**Original Caption**: a white surface with many yellow and indigo flowers

**System Generated Caption**: bouquet of colorful flowers in a small vase



**Original Caption**: a female tennis player shows her arm muscles

**System Generated Caption**: a male tennis player wearing white is playing tennis

# UNSUCCESSFUL CAPTIONS



**Original Caption**: a bunch of soccer players are playing a game

**System Generated Caption**: a baseball field with players and a crowd of spectators



**Original Caption**: a skateboarder riding their board in a skate park

**System Generated Caption**: diners at a cafe overlooking a sandy beach

# BLEU SCORE

BLEU: Bilingual Evaluation Understudy[4]

| BLEU-n | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 |
|--------|--------|--------|--------|--------|
| Score | 23.24 | 10.18 | 4.83 | 2.45 |

Fig: BLEU-n SCORE for the generated captions

# CONCLUSION

- System can successfully caption novel images for which representations are learnt.

- Since 1 caption per image was used for training, 5 caption per image would increase the accuracy manifold.

- System can generalize well for seen objects like train, pc, dog, snow etc.

# FUTURE WORKS

- Dense Image Captioning using better Object Identification Model, Inception V3[5].
- Description Generation for a given image.
- Generating t-SNE[6] representation for understanding vector space distance of similar captions. Similar captions will have smaller distance.
- Training the system on 200,000 images with 5 caption each.

21

# Thanks

# References

1. Karpathy, Andrej, and Li Fei-Fei. "Deep Visual-semantic Alignments for Generating Image Descriptions." 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2015). Web. 29 May 2016

2. Show and Tell: A Neural Image Caption Generator [https://github.com/karpathy/neuraltalk]

3. Microsoft COCO: Common Objects in Context - arXiv

4. BLEU: a Method for Automatic Evaluation of Machine Translation

5. Christian Szegedy et al. "Rethinking the Inception Architecture for Computer Vision"

6. Visualizing Data using t-SNE - Journal of Machine Learning Research