

# Image Captioning with Flickr30k Dataset

Meeting #3

Ye Kyaw Thu Lab Leader, LU Lab., Myanmar

# **Table of Content**

- Flickr30K Dataset Information
- Overview
- System Workflow & Coding Information

# Overview

- Image Team အတွက် ၂ရက်လောက် အချိန်ပေး စမ်းဖြစ်ခဲ့
- What is Image Captioning?
  - Automatically generate textual descriptions of images
- Key Components:
  - CNN (VGG16) for image features
  - RNN (GRU/LSTM) for sequence generation
  - Multilingual support (English/Myanmar)

#### Flickr30k Dataset

Estimated reading: 4 minutes • 8483 views



Visualization of the Flickr30k dataset in the Deep Lake UI

The Flickr30k dataset is a popular benchmark for sentence-based picture portrayal. The dataset is comprised of 31,783 images that capture people engaged in everyday activities and events. Each image has a descriptive caption. Flickr30k is used for understanding the visual media (image) that correspond to a linguistic expression (description of the image). This dataset is commonly used as a standard benchmark for sentence-based image descriptions.

Link: <a href="https://datasets.activeloop.ai/docs/ml/datasets/flickr30k-dataset/">https://datasets.activeloop.ai/docs/ml/datasets/flickr30k-dataset/</a>



- 1. a man prepares to enter the red building.
- 2. a man walking around the corner of a red building.
- 3. a man walks past a red building with a fake rocket attached to it.
- 4. a man walks under a building with a large rocket shaped sculpture.
- 5. a person walking by a red building with a jet on top of it .



- 1. a black dog playing with a purple toy in the snow.
- 2. a black dog runs through the snow carrying a blue toy.
- 3. a dog plays in the snow.
- 4. dog running with a purple toy in the snowy field.
- 5. the black and brown dog carries a purple toy in the snow.



- 1. a guy and a girl, both wearing white shirts and jeans, stand under a flowering tree.
- 2. a man and a woman are talking in a park
- 3. a man and woman standing underneath the tree are talking.
- 4. a man in a white shirt is standing in the grass showing something to a woman in a white shirt.
- 5. a young couple both wearing white shirts and blue jeans standing in a light misty rain

```
data
 captions. txt
  flickr30k images
      captions. txt
          - 1000092795. jpg
           10002456. jpg
           1000268201. jpg
           1000344755. jpg
           1000366164. jpg
           1000523639. jpg
          · 1000919630. jpg
          - 10010052. jpg
           1001465944. jpg
           1001545525. jpg
           1001573224. jpg
          - 1001633352. jpg
          - 1001773457. jpg
           1001896054. jpg
           100197432. jpg
          · 1003163366. jpg
           1003420127. jpg
```

- Folder structure of Flickr30k
- Label ဖိုင် နှစ်ဖိုင် ရှိတယ်
- Images ဖိုလ်ဒါအောက်မှာ
   ID နဲ့ သိမ်းထားတဲ့ image
   ဖိုင်တွေ ရှိလိမ့်မယ်

```
(base) ye@lst-hpc3090:~/intern3/img2txt/data$ head -n 20 ./captions.txt
                                 A child in a pink dress is climbing up a set of stairs in an entry way .
1000268201_693b08cb0e.jpg#0
1000268201_693b08cb0e.jpg#1
                                 A girl going into a wooden building .
                                 A little girl climbing into a wooden playhouse .
1000268201_693b08cb0e.jpg#2
1000268201_693b08cb0e.jpg#3
                                 A little girl climbing the stairs to her playhouse.
                                 A little girl in a pink dress going into a wooden cabin .
1000268201_693b08cb0e.jpg#4
1001773457_577c3a7d70.jpg#0
                                 A black dog and a spotted dog are fighting
1001773457_577c3a7d70.jpg#1
                                 A black dog and a tri-colored dog playing with each other on the road .
1001773457_577c3a7d70.jpg#2
                                 A black dog and a white dog with brown spots are staring at each other in the street
1001773457_577c3a7d70.jpg#3
                                 Two dogs of different breeds looking at each other on the road .
                                 Two dogs on pavement moving toward each other .
1001773457_577c3a7d70.jpg#4
1002674143_1b742ab4b8.jpg#0
                                 A little girl covered in paint sits in front of a painted rainbow with her hands in a
bowl .
1002674143_1b742ab4b8.jpg#1
                                 A little girl is sitting in front of a large painted rainbow .
                                 A small girl in the grass plays with fingerpaints in front of a white canvas with a ra
1002674143_1b742ab4b8.jpg#2
inbow on it .
1002674143_1b742ab4b8.jpg#3
                                 There is a girl with pigtails sitting in front of a rainbow painting.
                                 Young girl with pigtails painting outside in the grass .
A man lays on a bench while his dog sits by him .
1002674143_1b742ab4b8.jpg#4
1003163366_44323f5815.jpg#0
                                 A man lays on the bench to which a white dog is also tied .
1003163366_44323f5815.jpg#1
                                 a man sleeping on a bench outside with a white and black dog sitting next to him .
1003163366_44323f5815.jpg#2
                                 A shirtless man lies on a park bench with his dog .
1003163366_44323f5815.jpg#3
1003163366_44323f5815.jpg#4
                                 man laving on bench holding leash of dog sitting on ground
```

Fig. Flickr80k style

```
(base) ye@lst-hpc3090:~/intern3/img2txt/data/flickr30k_images$ head -n 20 captions.txt
image,caption
1000092795.\mathtt{jpg}, Two young guys with shaggy hair look at their hands while hanging out in the yard .
1000092795.jpg," Two young , White males are outside near many bushes ."
1000092795.jpg, Two men in green shirts are standing in a yard .
1000092795.jpg, A man in a blue shirt standing in a garden .
1000092795.jpg, Two friends enjoy time spent together .
10002456.jpg, Several men in hard hats are operating a giant pulley system .
10002456.jpg, Workers look down from up above on a piece of equipment.
10002456.jpg, Two men working on a machine wearing hard hats .
10002456.jpg, Four men on top of a tall structure .
10002456.jpg, Three men on a large rig .
1000268201.jpg, A child in a pink dress is climbing up a set of stairs in an entry way .
1000268201.jpg, A little girl in a pink dress going into a wooden cabin .
1000268201.jpg, A little girl climbing the stairs to her playhouse .
1000268201.jpg, A little girl climbing into a wooden playhouse
1000268201.jpg, A girl going into a wooden building .
1000344755.jpg, Someone in a blue shirt and hat is standing on stair and leaning against a window
1000344755.jpg, A man in a blue shirt is standing on a ladder cleaning a window .
1000344755.jpg, A man on a ladder cleans the window of a tall building .
1000344755.jpg, man in blue shirt and jeans on ladder cleaning windows
```

Fig. Flickr30k style

# System Workflow

## Visual Flow:

```
Image → VGG16 (Feature Extraction) →
GRU/LSTM (Caption Generation)
```

# Key Modules:

- Data loading
- Feature extraction
- Model training
- Evaluation

# Feature Extraction with VGG16

#### • Process:

- Resize images to 224x224
- Extract 4096-dim features from VGG16's penultimate layer

# Output:

- Saved as features.pkl for reuse
- Code
  - features = extract\_image\_features(image\_paths, vgg\_model)

# **Caption Tokenization**

#### Tokenizer

- $\circ$  Converts words to integers (e.g., "dog"  $\rightarrow$  42)
- Calculates vocab\_size and max\_length of captions

#### Code

o tokenizer.fit\_on\_texts(captions) # Vocab size = 10,000

# Train-Test Split

- Split Ratio: 90% train, 10% test
- Key Code:
  - o split = int(len(image\_ids) \* 0.9)
  - train\_keys = image\_ids[:split]
  - o test\_keys = image\_ids[split:]

## GRU vs. LSTM

- GRU (Gated Recurrent Unit):
  - Fewer parameters → Faster training
  - Single gate structure
- LSTM (Long Short-Term Memory):
  - More parameters → Better for long sequences
  - 3 gates (input, forget, output)
- Code Switch:
  - o if cell\_type == 'gru': se3 = GRU(units)
  - else: se3 = LSTM(units)

# Model Architecture

- Inputs:
  - Image features (4096-dim)
  - Caption sequences (padded to max\_length)
- Layers:
  - Embedding → Dropout → GRU/LSTM → Dense
- Loss:
  - Categorical cross-entropy

# **Training Process**

#### Parameters

- Batch size: 24
- o Epochs: 200
- Learning rate: 0.001, 0.002

## Code

generator = data\_generator(train\_keys, features, ...)

#### Metrics Used:

- BLEU (1-4), chrF++, CIDEr, ROUGE (1/2/L)
- Semantic similarity (TF-IDF + keyword overlap)

# Example Output:

- o BLEU-1: 0.2292
- ROUGE-L: 0.1530
- o CIDEr: 0.0003

- TF-IDF Cosine Similarity
  - Converts text to vectors (ignoring stopwords)
  - Measures contextual similarity via cosine\_similarity()
- Code
  - tfidf = vectorizer.fit\_transform(all\_texts)
  - sim = cosine\_similarity(tfidf[0:1], tfidf[1:2])[0][0]

# Keyword Overlap

 Exact word matches normalized by predicted caption length:

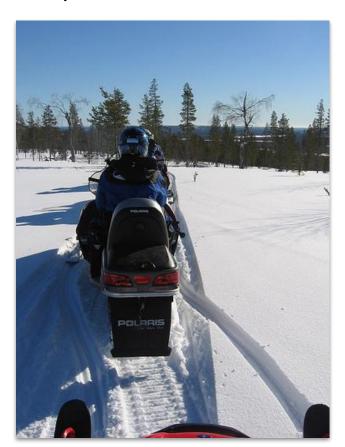
## Code

- len(pred\_words & ref\_words) / max(1, len(pred\_words))
- Composite Score
  - Weighted sum: 0.7 \* similarity + 0.3 \* keyword\_overlap

```
ef semantic_similarity(predicted, references):
  all texts = [' '.join(references), predicted]
   vectorizer = TfidfVectorizer(stop words='english')
      tfidf = vectorizer.fit transform(all texts)
   sim = cosine similarity(tfidf[0:1], tfidf[1:2])[0][0]
  pred_words = set(predicted.lower().split())
   ref words = set(' '. join(references).lower().split())
   keyword overlap = len(pred words & ref words) / max(1, len(pred words))
       'similarity': sim,
        keyword overlap': keyword overlap,
       composite score': 0.7*sim + 0.3*keyword overlap
```

# Results Analysis (Good Example 1)

• Image ID: 437527058



# Results Analysis (Good Example 1)

- Image ID: 437527058
- Actual: "a caravan of snowmobiles travels through snow"
- Predicted: "helmet and dog walking through snow"
- **Semantic Score:** 0.3424 (highest composite score)
- Why Good?: Partial keyword overlap ("snow") + contextually plausible

# Results Analysis (Good Example 2)

• Image ID: 441212506



# Results Analysis (Good Example 2)

- Image ID: 441212506
- Actual: "three dogs playing in a field"
- Predicted: "dog with frisbee rolling in grass"
- Semantic Score: 0.3059
- Why Good?: Captured "dog" and outdoor activity

# Results Analysis (Good Example 3)

• **Image ID**: 440737340



# Results Analysis (Good Example 3)

- Image ID: 440737340
- Actual: "masked man carrying a box"
- Predicted: "tattoos sitting on a bench"
- Semantic Score: 0.0321 (lowest)
- Why Good?: Complete mismatch

# CLI Usage

- Commands:
- # Train
  - python image\_captioning.py --train --epochs 200--cell\_type gru
- # Predict
  - python image\_captioning.py --predict img.jpg

# CLI Usage

```
Image Captioning with Flickr30k Dataset
optional arguments:
 -h, --help
                        show this help message and exit
 --data dir DATA DIR
                        Directory to store dataset (default: ./data)
 --model dir MODEL DIR
                        Directory to save/load models (default: ./models)
 --epochs EPOCHS
                        Number of training epochs (default: 15)
 --batch size BATCH SIZE
                        Training batch size (default: 64)
 --train
                       Train the model
 --evaluate
                        Evaluate on test set
 --predict PREDICT
                        Path to single image for prediction
 --language {english, myanmar}
                        Caption language (default: english)
 --model name MODEL NAME
                        Model filename (default: best model. h5)
 --skip download
                       Skip dataset download (for prediction mode)
 --skip feature extraction
                        Skip feature extraction if features file exists
 -- 1stm units LSTM UNITS
                        Number of units in LSTM layer (default: 256)
 --dropout rate DROPOUT RATE
                        Dropout rate (default: 0.4)
 --learning rate LEARNING RATE
                        Learning rate (default: 0.001)
 --early stopping EARLY STOPPING
                        Patience for early stopping (default: None)
 -- feature size FEATURE SIZE
                        Size of image features (default: 4096)
 --cell type {lstm, gru}
                        RNN cell type (default: gru)
```

## To Do

## Transformer-based models

o အချိန်ရတဲ့အခါ လက်ရှိ code ကိုပဲ အခြေခံပြီး Transformer architecture နဲ့ ဖြည့်ရေးတာလုပ်နိုင်

# Better Myanmar language support

o လက်ရှိ ဒေတာကိုလည်း မြန်မာလို ဘာသာပြန်လိုက်ပြီး အဲဒါကို experiment အနေနဲ့ လုပ်ကြည့်တာမျိုး