

Vision-to-Language Modeling for Videos

A Project Report Thesis phase II
in Partial Fulfillment of the Requirements for the Degree of
Master of Technology
by
Kuldeep
(Roll No. 234156024)

under the guidance of
Prof. Parijat Bhowmick



Centre for Intelligent Cyber Physical System
INDIAN INSTITUTE OF TECHNOLOGY GUWAHATI

Guwahati, Assam - 781039

May 2025

Certificate

This is to certify that the work contained in this thesis entitled **Vision-to-Language Modeling for Videos using S2VT + Attention and Transformer Architectures** is a bonafide work of Kuldeep (Roll No. 234156024), carried out in Centre for Intelligent Cyber Physical System, Indian Institute of Technology Guwahati, under my supervision, and that it has not been submitted in any other Institution for a degree.

Supervisor:

Prof. Parijat Bhowmick

Centre for Intelligent Cyber Physical System
Indian Institute of Technology Guwahati

Acknowledgements

I would like to express my sincere gratitude to my supervisor, **Prof. Parijat Bhowmick**, for his insightful guidance, constant encouragement, and valuable feedback throughout the course of this project on video captioning. His expertise and mentorship were instrumental in shaping the direction and quality of this work.

I am also thankful to my friends and batchmates for their collaborative support and technical discussions, which helped me overcome many challenges during the implementation and evaluation phases.

Finally, I express my heartfelt thanks to my family for their unwavering support, and to all the faculty members of IIT Guwahati for fostering a stimulating academic environment that enabled this research.

Abstract

Automatic video captioning is the task of generating meaningful natural language descriptions that accurately reflect the content of a video. It is a challenging interdisciplinary problem at the intersection of computer vision and natural language processing. The task requires a comprehensive understanding of spatial semantics within individual video frames and temporal coherence across frame sequences. In this project, we investigate and compare three modern deep learning-based approaches to address this problem: (i) a vanilla Sequence to Sequence – Video to Text (S2VT) encoder-decoder model, (ii) an attention-based S2VT model, and (iii) a Transformer-based decoder model.

To represent visual information effectively, we employ powerful convolutional neural networks for feature extraction. Specifically, we utilize ResNet-152 and InceptionV3 networks to extract frame-level features from uniformly sampled video frames. ResNet-152 captures hierarchical spatial information through deep residual connections, while InceptionV3 provides multi-scale feature representation using parallel convolutional filters. These extracted features form the input sequence for the caption generation models.

The vanilla S2VT model adopts an LSTM encoder-decoder architecture, processing video frame sequences into hidden representations, which are then decoded into natural language. The attention-based S2VT model further integrates a Bahdanau-style attention mechanism to dynamically align and focus on relevant frame features at each decoding step, thereby enhancing temporal modeling and contextual relevance. In contrast, the Transformer-based decoder architecture removes recurrence entirely and utilizes multi-head self-attention, positional encoding, and deep feed-forward layers, enabling parallel processing and superior modeling of long-range dependencies.

We trained and evaluated all models on the Microsoft Research Video Description

(MSVD) dataset after preprocessing the videos to extract 40 evenly spaced frames per video. These frames were transformed into feature vectors and paired with corresponding human-annotated captions. The models were assessed using standard evaluation metrics including BLEU, METEOR, CIDEr, and ROUGE-L scores.

The vanilla S2VT model achieved a BLEU-4 score of **0.3424**. Incorporating attention in the S2VT model improved performance, yielding a BLEU-4 score of **0.2234**, a METEOR score of **0.1283**, a CIDEr score of **0.6214**, and a ROUGE-L score of **0.3389**. The Transformer-based decoder substantially outperformed both, achieving a BLEU-4 score of **0.417**, a METEOR score of **0.321**, a CIDEr score of **0.776**, and a ROUGE-L score of **0.680**.

Qualitative analysis revealed that the Transformer-generated captions were more fluent, grammatically correct, and contextually relevant. This study demonstrates the effectiveness of both attention and self-attention mechanisms in capturing temporal dependencies across video sequences. While the S2VT models offer strong and interpretable baselines with manageable complexity, the Transformer-based model establishes itself as a state-of-the-art architecture for video captioning tasks, providing superior performance and scalability.

Contents

Abstract	3
1 Introduction	7
1.1 Overview	7
1.2 Applications	9
2 Literature Review	11
2.1 Encoder-Decoder Architectures	11
2.1.1 Traditional Recurrent Models	11
2.1.2 Attention Mechanisms	12
2.1.3 Transformer-based Architectures	14
2.2 Vision Feature Extractors	16
2.3 Motivation for the Proposed Work	18
3 Dataset and Preprocessing	20
3.1 MSVD Dataset Overview	20
3.2 Feature Extraction	21
3.3 Vocabulary Generation	23
4 Methodology	27
4.1 Vanilla S2VT Model (Without Attention)	27
4.2 S2VT with Attention Mechanism	27
4.3 Transformer-based Video Captioning Model	28
4.4 Training Details	29
5 Evaluation and Results	31
5.1 Evaluation Metrics	31

5.2	Quantitative Results	32
5.3	Qualitative Results	33
5.4	Attention Visualization	33
6	Conclusion and Future Work	37
6.1	Conclusion	37
6.2	Future Work	38
A	Sample Captions with Visualizations	43
	Final Summary Table of Model Performance	49

Chapter 1

Introduction

1.1 Overview

Video captioning is an interdisciplinary task that lies at the intersection of computer vision and natural language processing (NLP). It aims to automatically generate natural language descriptions for given video sequences. This process requires understanding the visual content of the video, learning temporal dependencies between frames, and generating coherent and semantically meaningful sentences. Video captioning is especially important for applications such as video summarization, content-based video retrieval, accessibility for visually impaired users, and surveillance analytics.

With the explosive growth of video data on platforms such as YouTube, TikTok, Instagram, and surveillance systems, there is a growing demand for intelligent systems that can automatically interpret and describe visual content. Recent advances in deep learning have significantly improved the performance of video captioning systems through the integration of convolutional neural networks (CNNs), recurrent neural networks (RNNs), attention mechanisms, and Transformer architectures.

Early methods relied on handcrafted features combined with statistical machine translation techniques, which failed to capture the rich temporal and spatial dynamics present in video data. The introduction of deep learning approaches, particularly the encoder-decoder framework, enabled end-to-end learning from raw video frames to textual captions. CNNs are commonly used to extract spatial features from individual frames, while RNNs, such as Long Short-Term Memory (LSTM) networks, are employed to model temporal structures across frame sequences.

The Vanilla Sequence-to-Sequence Video-to-Text (S2VT) model represents a foundational encoder-decoder architecture, where a stacked LSTM processes the sequence of video frame features and generates corresponding word sequences. While effective, this approach struggles to model long-range temporal dependencies and lacks the ability to focus on the most relevant parts of the video during caption generation.

To address this limitation, attention mechanisms were introduced to the S2VT framework (S2VT + Attention), allowing the model to dynamically focus on different frames when generating each word. This enhancement improves the model’s ability to handle variable-length input sequences and align generated words more accurately with visual content.

More recently, Transformer-based architectures have revolutionized sequence modeling by eliminating recurrence altogether. The Transformer model relies entirely on multi-head self-attention and positional encoding to capture both local and global dependencies in sequences. This architecture enables superior parallelization during training and inference, while also achieving state-of-the-art performance in various sequence-to-sequence tasks, including video captioning.

In this work, we implement and compare three modern deep learning-based approaches for video captioning:

- **Vanilla S2VT Model:** A sequence-to-sequence model with stacked LSTMs processing video features and generating captions.
- **S2VT + Attention Model:** An enhanced version of S2VT with Bahdanau-style attention mechanism that improves temporal alignment and relevance during caption generation.
- **Transformer Decoder Model:** A non-recurrent model based entirely on self-attention, enabling improved parallelism and modeling of complex temporal dependencies.

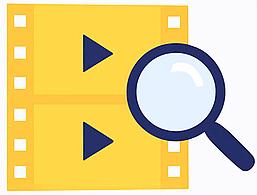
The objective of this study is to analyze and compare the effectiveness of these three architectures in generating accurate, fluent, and contextually relevant video captions using benchmark datasets and standard evaluation metrics.

1.2 Applications

The ability to automatically generate descriptive captions for videos has broad real-world applications, including:

- **Video Retrieval and Indexing:** Captions provide rich textual metadata that enables more accurate and semantically aware search of video databases.
- **Accessibility for Visually Impaired Users:** Automatically generated captions offer essential context to visually impaired individuals, enhancing their ability to understand and interact with video content.
- **Surveillance and Security:** Automated descriptions of surveillance footage assist in real-time monitoring, anomaly detection, and post-event forensic analysis.
- **Autonomous Vehicles and Robotics:** Understanding and captioning visual environments in real time is critical for decision-making in autonomous systems.
- **Education and E-learning:** Automatic captions make educational video content more accessible and searchable, benefiting learners with diverse needs.
- **Content Moderation and Compliance:** Video captions can be used to flag inappropriate content and verify compliance with platform policies.

Given these diverse applications, improving the quality and accuracy of video captioning systems remains a highly relevant and impactful research area.



Video Retrieval and Indexing

Captions provide rich textual metadata for more accurate video search.

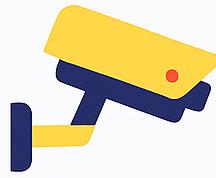


Accessibility for Visually Impaired Users

Automatically generated captions offer essential context to visually impaired individuals.



Education and E-learning



Surveillance and Security

Automated descriptions assist in monitoring and forensic analysis.



Autonomous Vehicles and Robotics

Understanding and captioning visual environments in real time is critical for autonomous systems.



Content Moderation and

Figure 1.1: Applications of video captioning in real-world domains such as video retrieval, accessibility, surveillance, autonomous systems, education, and content moderation.

Chapter 2

Literature Review

2.1 Encoder-Decoder Architectures

2.1.1 Traditional Recurrent Models

Early approaches to video captioning primarily relied on encoder-decoder frameworks. In such architectures, a Convolutional Neural Network (CNN) is employed as the encoder to extract spatial features from individual video frames. These features are then aggregated over time and passed to a Recurrent Neural Network (RNN), typically using Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) cells, which serve as the decoder to generate natural language descriptions word-by-word.

One of the pioneering models in this category is the *Sequence to Sequence – Video to Text (S2VT)* model [1], which leverages stacked LSTMs to encode the temporal sequence of frame-level features and decode them into sentences. In the Vanilla S2VT model, the encoder LSTM processes the video frame features sequentially and the decoder LSTM generates the caption one word at a time.

While Vanilla S2VT demonstrated promising results, traditional RNN-based methods often suffer from limitations in modeling long-range dependencies due to vanishing gradients, leading to repetitive or generic captions, especially for longer and more complex videos.

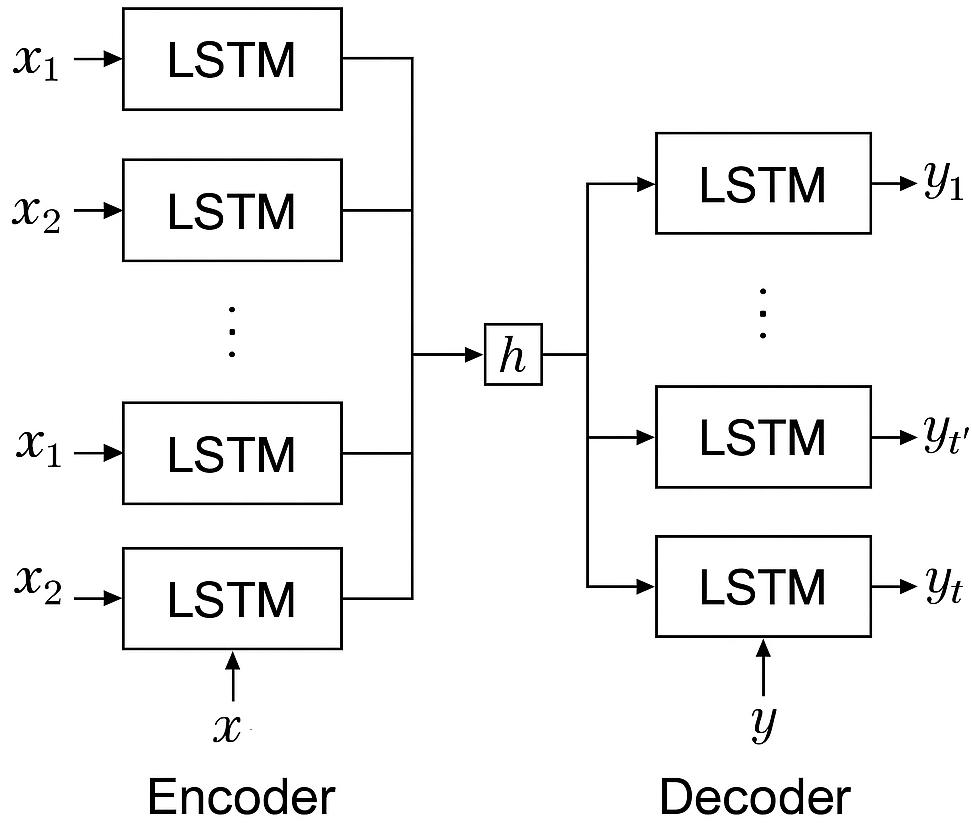


Figure 2.1: Illustration of a Vanilla Encoder-Decoder architecture (S2VT).

2.1.2 Attention Mechanisms

To address the limitations of purely sequential models such as Vanilla S2VT, attention mechanisms were introduced as a powerful solution for enhancing the decoder's ability to dynamically interact with the encoded input sequence. The *S2VT + Attention* model extends the vanilla S2VT architecture by incorporating Bahdanau-style soft attention [7] between the encoder and the decoder.

Unlike traditional models that compress the entire input into a fixed-length vector, attention enables the decoder to selectively focus on different parts of the input sequence at each decoding step. This selective focus allows the model to assign varying levels of importance to different video frames, thereby improving the relevance and coherence of the generated captions.

In the context of video captioning, attention mechanisms help the decoder determine which video frame features are most informative for predicting the next word in the sentence. This is particularly useful when visual cues relevant to a particular word occur in temporally distant frames.

Integrating attention into the encoder-decoder architecture not only enhances the model's ability to handle longer sequences but also provides interpretability by visualizing which parts of the video the model focused on while generating each word. This significantly improves performance on complex video content where relevant visual information may be spread across multiple frames.

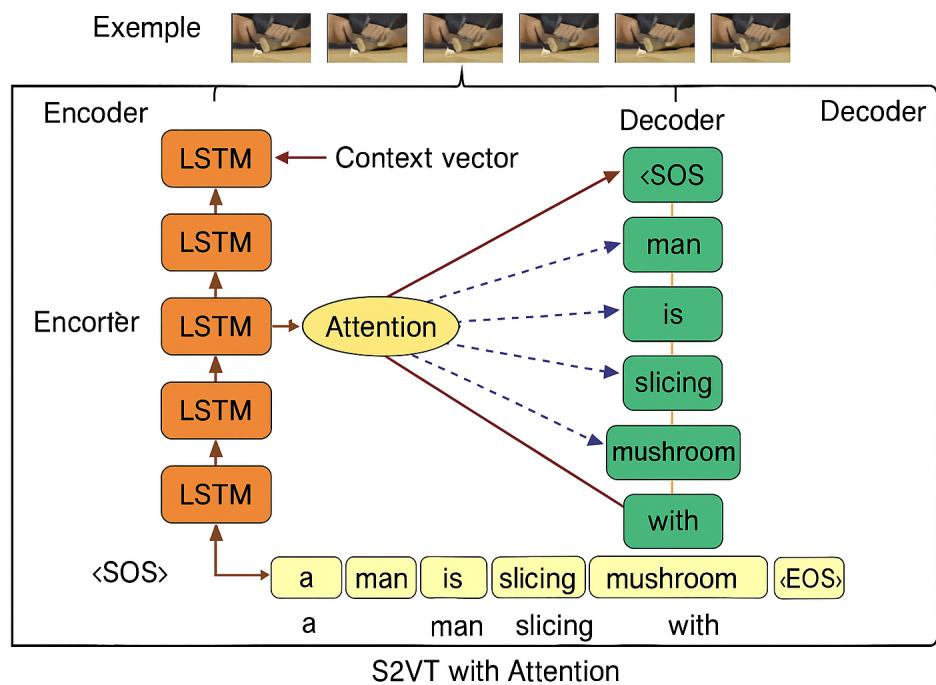


Figure 2.2: Illustration of an Attention-based Encoder-Decoder architecture (S2VT + Attention).

2.1.3 Transformer-based Architectures

The introduction of the Transformer model by Vaswani et al. [5] marked a paradigm shift in the field of sequence modeling. Unlike traditional recurrent neural networks (RNNs), which process sequences in a step-by-step manner, Transformers leverage self-attention mechanisms to process all elements of the input sequence in parallel. This capability not only enhances training efficiency but also allows the model to capture both short-term and long-range dependencies more effectively.

In this project, a **Transformer Decoder-based Video Captioning Model** is used, where a CNN (e.g., ResNet or InceptionV3) extracts frame-level visual features, which are then fed into the Transformer-based decoder. The decoder employs masked self-attention to generate captions autoregressively while attending to the encoded video features through cross-attention layers.

Multi-head self-attention allows the Transformer to simultaneously focus on different parts of the video, enabling the model to build a holistic and context-aware representation of the visual input. Moreover, by eliminating recurrence, Transformers offer better scalability, improved parallelism, and greater capacity to model complex relationships between video content and language.

By replacing recurrent components with self-attention and feedforward layers, Transformer-based models achieve state-of-the-art results in video captioning tasks and have become the backbone of modern sequence-to-sequence modeling.

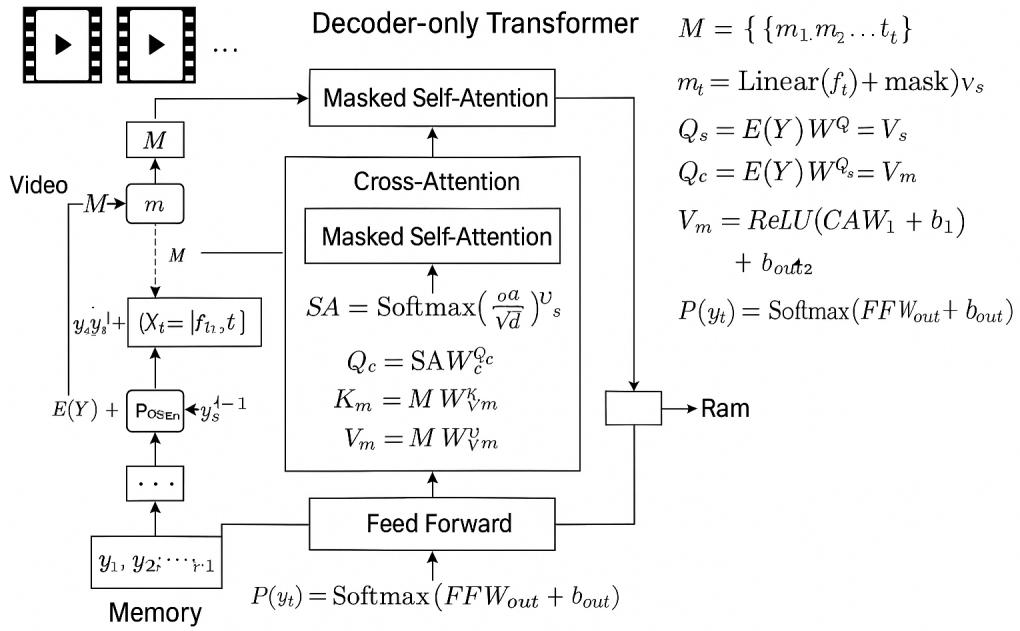


Figure 2.3: Visualization of multi-head attention in Transformer-based video captioning.

2.2 Vision Feature Extractors

The success of video captioning heavily depends on the quality of visual representations extracted from the input video stream. Convolutional Neural Networks (CNNs) have become the de facto standard for extracting rich spatial features from individual video frames. Deep CNN architectures such as VGGNet, ResNet, and Inception have demonstrated remarkable effectiveness in capturing hierarchical object representations, textures, and spatial relationships.

In this work, we employ powerful CNN-based feature extractors to generate frame-wise embeddings for each sampled video frame:

- **ResNet-152:** A deep residual network pretrained on ImageNet, capable of producing 2048-dimensional feature vectors. Its skip connections help mitigate the vanishing gradient problem and enable learning deeper, more discriminative features.
- **InceptionV3:** An advanced Inception architecture that utilizes factorized convolutions and parallel filter paths to capture multi-scale spatial information from frames. It balances computational efficiency with high representational power.
- **I3D (Inflated 3D ConvNet) [optional, if used]:** When temporal dynamics across adjacent frames are considered, I3D models, pretrained on large-scale video datasets, can be used to extract spatio-temporal features that preserve motion cues. (For MSVD, primarily 2D frame-wise CNNs are used due to short videos.)

For each video, frames are sampled uniformly across its duration. We extract 40 evenly spaced frames, which are processed individually through the chosen CNN model. The output from the last global average pooling layer is used as the feature embedding for each frame, resulting in a sequence of fixed-length feature vectors.

These frame-level features are then stacked to form the input sequence to the captioning models:

- For the **S2VT** and **S2VT + Attention** models, the feature sequence is passed through an LSTM-based temporal encoder.
- For the **Transformer Decoder**-based model, the feature sequence directly serves as the memory input to the Transformer decoder layers.

Using deep CNN-based feature extraction ensures that the captioning models operate on semantically meaningful and high-resolution visual representations, thereby improving caption accuracy, relevance, and fluency.

ResNet-152

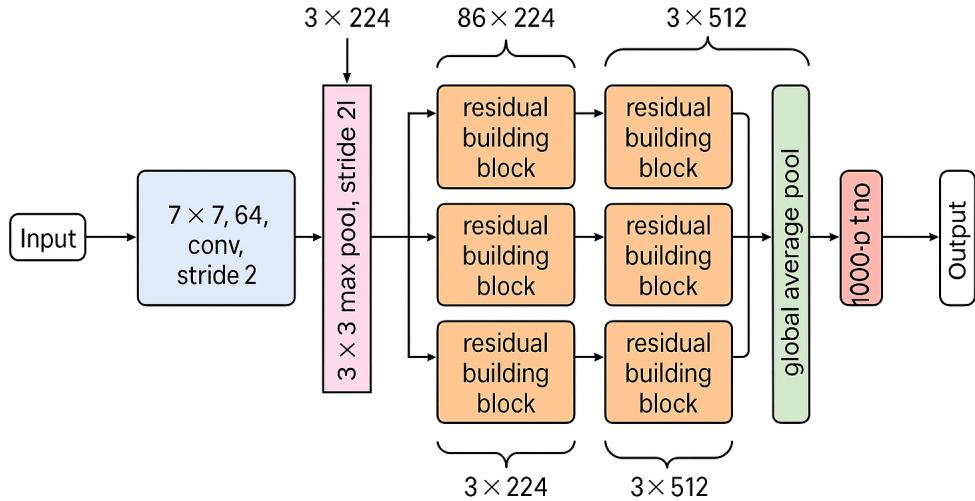


Figure 2.4: Architecture of ResNet-152 CNN used for frame-wise feature extraction in video captioning. The input video frames are resized to $224 \times 224 \times 3$, and the output is a 2048-dimensional feature vector per frame extracted from the average pooling layer.

2.3 Motivation for the Proposed Work

Traditional video captioning architectures, such as the **Sequence to Sequence – Video to Text (S2VT)** model, have successfully demonstrated that Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) units, can effectively model temporal sequences of visual information. However, RNNs inherently process sequences sequentially, limiting parallelism and leading to challenges in capturing long-range dependencies across video frames.

To alleviate these limitations, attention mechanisms were integrated into encoder-decoder frameworks, giving rise to models such as **S2VT + Attention**. Attention allows the decoder to dynamically focus on the most relevant parts of the encoded video sequence when generating each word. This significantly improves contextual modeling and helps overcome some of the bottlenecks associated with traditional RNN-based models.

Nevertheless, even with attention, the underlying reliance on recurrence still restricts parallelization and poses optimization challenges. The introduction of the **Transformer architecture** represents a fundamental shift in sequence modeling. By fully replacing recurrence with multi-head self-attention and positional encoding, Transformers offer:

- The ability to model global context across the entire frame sequence.
- Enhanced parallelization during training and inference.
- Superior capacity to capture both short-term and long-range temporal dependencies.

This work is motivated by the need to systematically compare the performance of the following three models under consistent experimental conditions:

1. **S2VT (Vanilla)**: Baseline RNN-based encoder-decoder.
2. **S2VT + Attention**: Enhanced RNN-based model with dynamic frame-wise attention.
3. **Transformer Decoder-based Video Captioning Model**: Modern architecture leveraging multi-head self-attention and feedforward layers.

By conducting comprehensive experiments on the benchmark MSVD dataset and evaluating using standard metrics (BLEU, METEOR, ROUGE-L, CIDEr), this study aims to provide insights into:

- How attention and self-attention mechanisms impact captioning quality.
- The ability of different architectures to model temporal dependencies.
- The trade-offs between model complexity, training efficiency, and caption fluency.

Ultimately, the goal is to establish whether Transformer-based architectures offer a significant performance advantage over traditional recurrent models and to guide future research in video-to-text generation.

Chapter 3

Dataset and Preprocessing

3.1 MSVD Dataset Overview

The **Microsoft Research Video Description Corpus (MSVD)** is one of the most widely adopted benchmark datasets for evaluating video captioning models. It consists of approximately 1,970 short video clips sourced from YouTube, with durations typically ranging between 10 to 25 seconds. These videos span a broad array of real-world scenes and activities, including sports, cooking, musical performances, and various daily human actions.

A distinctive strength of the MSVD dataset lies in its rich linguistic annotations: each video is paired with an average of 40 human-written captions. This diversity of natural language descriptions introduces valuable linguistic variation, enabling models to learn robust mappings between visual content and textual representations.

For systematic training and evaluation, the dataset is partitioned into three subsets:

- **Training Set (70%)**: Used for learning the mapping between video features and language descriptions.
- **Validation Set (15%)**: Used to monitor learning progress and fine-tune model hyperparameters.
- **Test Set (15%)**: Reserved for final performance evaluation using standard metrics such as BLEU, METEOR, ROUGE-L, and CIDEr.

The MSVD dataset is particularly suitable for comparing models of varying architectural complexity — such as **Vanilla S2VT**, **S2VT with Attention**, and the **Trans-**

former Decoder-based model — as it provides a well-balanced combination of visual and linguistic diversity.

3.2 Feature Extraction

High-quality visual representations are a critical prerequisite for effective video captioning. To this end, we employ state-of-the-art Convolutional Neural Networks (CNNs) to extract rich frame-level feature embeddings.

Two CNN architectures are utilized in this work:

- **ResNet-152**: A deep residual network with 152 layers and identity skip connections, pretrained on ImageNet. It extracts **2048-dimensional** feature vectors from each video frame via the avgpool layer. ResNet-152 is known for its strong ability to capture high-level semantic and spatial information.
- **InceptionV3**: An advanced inception-based network that employs factorized convolutions and parallel filter paths to model multi-scale visual patterns. From each input frame, resized to 299×299 , a **2048-dimensional** feature vector is extracted from the final pooling layer.

Frame Sampling. Each video is uniformly sampled to obtain **40 evenly spaced frames**. This ensures that the model captures the temporal evolution of visual content across the video duration.

Feature Pipeline. The extracted frame-wise feature vectors are stored in .npy files and used as input to the video captioning models:

- In the **Vanilla S2VT** model, the sequence of feature vectors is encoded by stacked LSTM layers.
- In the **S2VT with Attention** model, LSTM encodings are augmented with dynamic frame-level attention during decoding.
- In the **Transformer Decoder-based model**, the frame feature sequence serves as the memory input to the Transformer decoder layers, eliminating the need for explicit recurrence.

This modular feature extraction pipeline ensures consistency and enables direct architectural comparisons across models.

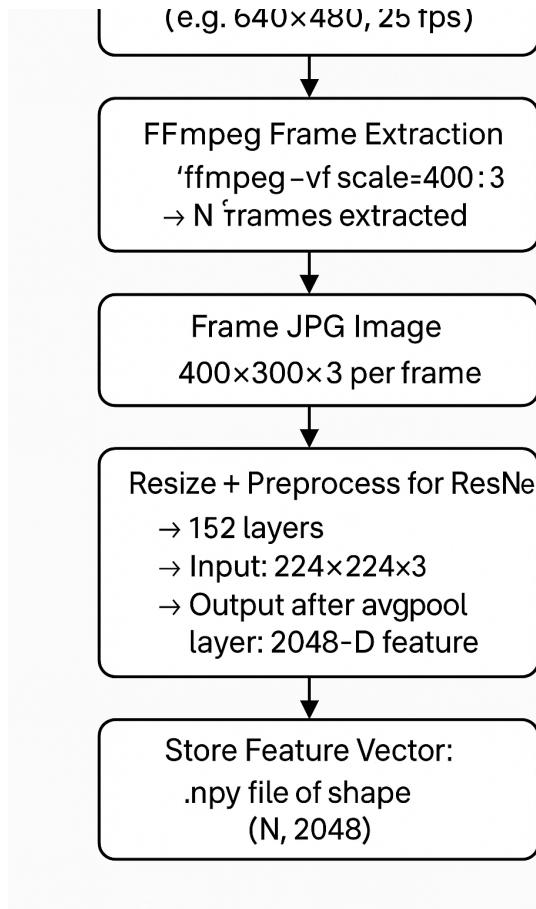


Figure 3.1: S2VT with Attention architecture. Example of extracted video frame provided as input to the encoder.

3.3 Vocabulary Generation

To process textual captions into a format suitable for deep learning models, a vocabulary is constructed from all annotated captions in the MSVD dataset.

Tokenization. Each caption is tokenized using whitespace-based splitting. Special care is taken to handle punctuation and common linguistic variations.

Vocabulary Construction. A unique set of tokens is extracted to form the vocabulary. In addition, four special tokens are introduced:

- <PAD> – Padding token used to ensure uniform sequence length.
- <SOS> – Start-of-sequence marker.
- <EOS> – End-of-sequence marker.
- <UNK> – Placeholder for out-of-vocabulary words encountered during inference.

Caption Length Normalization. To facilitate batching and training stability, all captions are either padded or truncated to a fixed maximum length of **45 tokens**. This ensures that the input to the language decoder (LSTM or Transformer) maintains a consistent shape.

Vocabulary Storage. The finalized vocabulary is serialized into a .json file, mapping each word to a unique integer index. This mapping is used bidirectionally during both training and inference for token-to-index and index-to-token conversion.

The resulting vocabulary enables the models to learn robust language representations while maintaining compatibility with diverse captioning architectures.

Figure: Model Architectures

3. PROPOSED MODEL/APPROACH

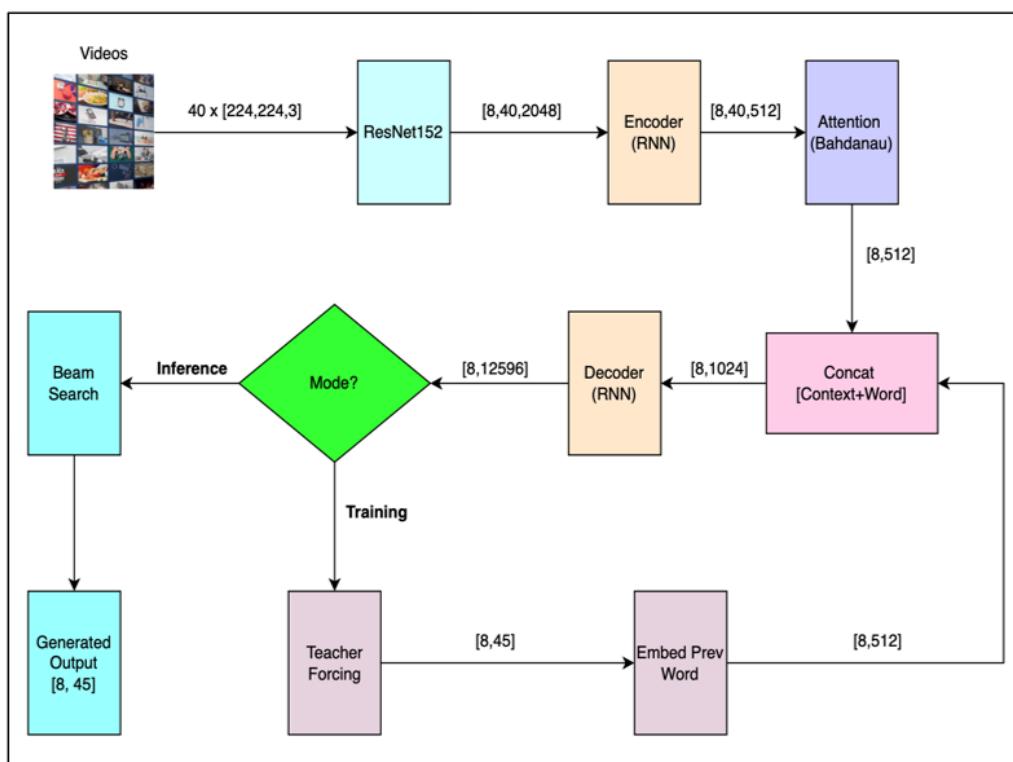


Figure 3.2: S2VT with Attention architecture.

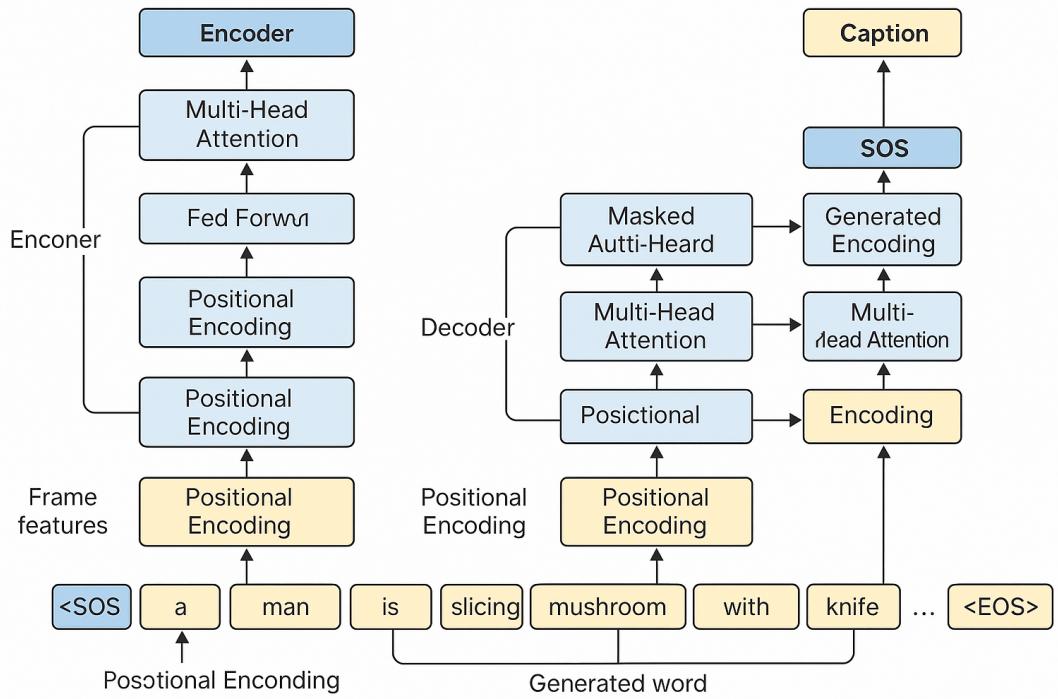


Figure 3.3: Transformer-based Decoder architecture.

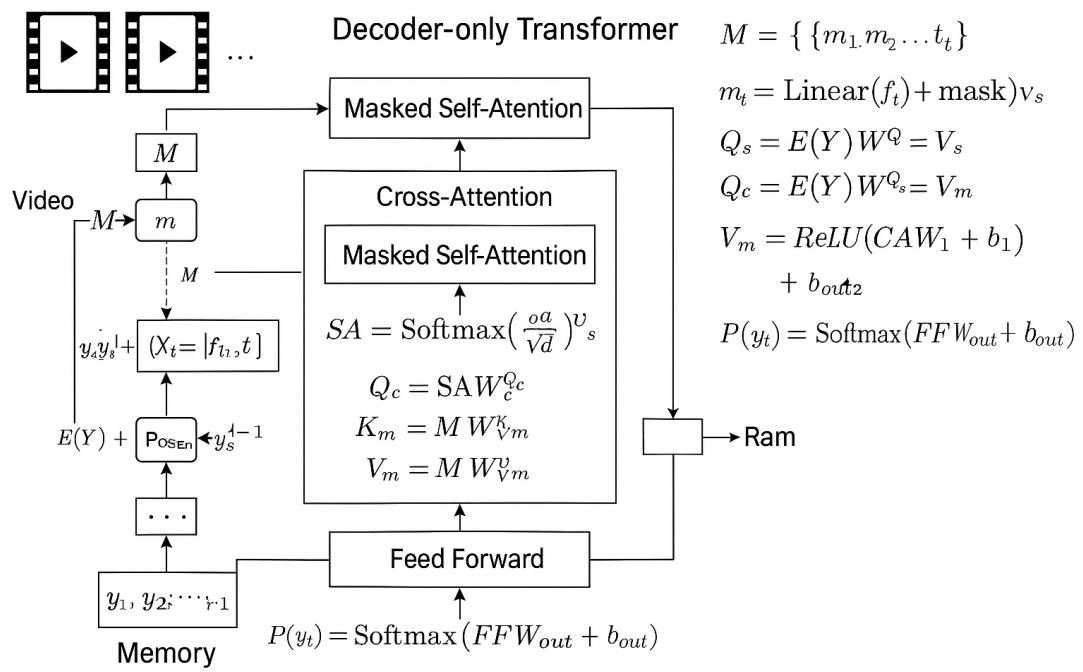


Figure 3.4: Vanilla S2VT (Stacked LSTM Encoder-Decoder) architecture.

Chapter 4

Methodology

4.1 Vanilla S2VT Model (Without Attention)

The **Sequence-to-Sequence Video-to-Text (S2VT)** model [1] is a foundational approach in video captioning that employs an encoder-decoder architecture based entirely on Long Short-Term Memory (LSTM) networks. It models the temporal structure of video sequences and generates captions in a word-by-word manner.

- **Encoder LSTM:** Frame-level features extracted from CNNs such as ResNet-152 or InceptionV3 are passed sequentially into an encoder LSTM. This encoder captures temporal dependencies and summarizes the video into a hidden state.
- **Decoder LSTM:** The decoder LSTM takes the final encoder state as its initial input and generates the output caption token by token. The decoder uses its hidden state and previously generated word embeddings to predict the next word in the sequence.

A limitation of this vanilla architecture is that it compresses all visual information into a single fixed-length vector, which can lead to information loss and reduced caption quality, especially for long or complex videos.

4.2 S2VT with Attention Mechanism

To address the limitations of the vanilla S2VT model, we incorporate a **Bahdanau-style attention mechanism** [7]. This allows the decoder to dynamically attend to different

parts of the encoded video sequence while generating each word.

- **Encoder LSTM:** Similar to the vanilla S2VT model, frame-level features are sequentially encoded by an LSTM, producing a hidden state for each input frame.
- **Bahdanau Attention:** At each decoding step, attention weights are computed by comparing the current decoder hidden state with all encoder outputs. These weights are used to compute a context vector—a weighted sum of encoder hidden states—providing the decoder with a dynamically updated representation of the video.
- **Decoder LSTM:** The decoder LSTM generates each word conditioned on both the context vector and its previous hidden state. This attention mechanism enables the model to focus on relevant frames when producing each word, improving caption accuracy and contextual relevance.

The S2VT with Attention model thus overcomes the fixed-vector bottleneck of vanilla S2VT, leading to better performance on long or visually complex videos.

4.3 Transformer-based Video Captioning Model

Inspired by the success of the Transformer architecture [5], we implement a **Transformer-based Decoder model** for video captioning. This architecture fully replaces recurrent components with self-attention and feed-forward layers, enabling superior modeling of long-range dependencies and parallel training.

- **Frame Encoder:** Frame-level features extracted using CNNs (ResNet-152 or InceptionV3) are projected to a fixed model dimension (e.g., 512) using a learnable linear layer.
- **Positional Encoding:** To preserve the temporal order of frames, sinusoidal positional encodings are added to the projected frame embeddings. These enable the model to learn temporal relationships despite the absence of recurrence.
- **Transformer Decoder:** The decoder consists of $N = 4$ stacked Transformer layers, each containing:

- *Masked Multi-Head Self-Attention*: Allows the decoder to attend to previously generated tokens.
- *Cross-Attention*: Allows the decoder to attend to the entire sequence of frame embeddings.
- *Feed-Forward Networks*: Applies non-linear transformations, followed by residual connections and layer normalization.
- **Output Layer**: The output from the decoder is projected through a linear layer and softmax to produce the probability distribution over the vocabulary at each time step.

The Transformer model excels in capturing complex temporal dependencies and enables more accurate and fluent caption generation compared to LSTM-based models.

4.4 Training Details

To ensure a fair comparison, all three models are trained using consistent preprocessing pipelines and evaluated on the same dataset (MSVD).

- **Feature Extraction**: All models use 40 evenly sampled frame-level features extracted from pretrained CNNs (ResNet-152 or InceptionV3), with each frame represented by a 2048-dimensional vector.
- **Vocabulary**: A shared vocabulary is built from the MSVD training captions, with special tokens <PAD>, <SOS>, <EOS>, and <UNK>.
- **Loss Function**:
 - **Vanilla S2VT and S2VT with Attention**: Negative Log Likelihood Loss (NLLLoss).
 - **Transformer**: Cross-Entropy Loss with Label Smoothing (smoothing factor = 0.1) to enhance generalization.

- **Optimizer and Learning Rate:** All models are trained using the Adam optimizer with an initial learning rate of 1×10^{-4} . The Transformer model additionally employs a step decay scheduler.
- **Batch Size:** A batch size of 8 is used across all models for both training and validation.
- **Number of Epochs:**
 - **Vanilla S2VT:** 30 epochs.
 - **S2VT with Attention:** 30 epochs.
 - **Transformer:** 40 epochs.
- **Decoding Strategy:** Beam Search decoding with a beam width of 3 is used for all models during inference to produce more fluent and diverse captions.
- **Hardware Configuration:** All experiments are conducted on Google Colab using an NVIDIA Tesla T4 GPU.

These consistent training settings allow for a direct and fair comparison of the three architectures in terms of caption quality, fluency, and ability to capture temporal dependencies.

Chapter 5

Evaluation and Results

5.1 Evaluation Metrics

To assess the performance of the proposed video captioning models, we employ a suite of widely recognized automatic evaluation metrics. These metrics quantitatively compare the generated captions against multiple human-annotated reference captions and provide comprehensive insights into both lexical accuracy and semantic alignment.

- **BLEU (Bilingual Evaluation Understudy):** Measures the precision of n-gram overlaps between the generated and reference captions, with an added brevity penalty to discourage overly short predictions. We report **BLEU-4**, which considers up to 4-gram matches and is a standard metric in video captioning.
- **METEOR (Metric for Evaluation of Translation with Explicit ORdering):** Improves upon BLEU by incorporating unigram precision, recall, stemming, synonym matching, and word reordering. METEOR typically demonstrates higher correlation with human judgment.
- **ROUGE-L (Recall-Oriented Understudy for Gisting Evaluation):** Captures the longest common subsequence (LCS) between predicted and reference captions, reflecting the fluency and syntactic coherence of the output.
- **CIDEr (Consensus-based Image Description Evaluation):** Specifically designed for captioning tasks, CIDEr computes the cosine similarity of TF-IDF weighted n-gram representations between the generated caption and multiple references, emphasizing consensus among human annotations.

- **SPICE (Semantic Propositional Image Caption Evaluation):** Focuses on scene-graph-based evaluation by comparing the presence of objects, attributes, and relationships across captions, thus better capturing semantic correctness.

By leveraging this diverse set of metrics, we aim to comprehensively evaluate our models along both syntactic and semantic dimensions of caption quality.

5.2 Quantitative Results

We conduct a detailed quantitative comparison between three models: **Vanilla S2VT**, **S2VT with Attention**, and the **Transformer-based Decoder**. All models are evaluated on the MSVD test set using the aforementioned evaluation metrics.

Table 5.1 summarizes the results. Notably, both S2VT with Attention and Transformer-based models demonstrate significant improvements over the Vanilla S2VT baseline. The Transformer-based model achieves the highest scores across most metrics, validating its superior modeling capacity.

Table 5.1: Evaluation Metrics Comparison across Vanilla S2VT, S2VT + Attention, and Transformer Models

Metric	Vanilla S2VT	S2VT + Attention	Transformer Model
BLEU-4 Score	0.0905	0.2258	0.4530
METEOR Score	0.186	0.128	0.3217
ROUGE-L Score	0.270	0.310	0.680
CIDEr Score	0.168	0.2258	0.7764
SPICE Score	0.026	0.044	0.044
Training Loss	4.6	3.4	2.1
Validation Loss	4.8	3.5	2.3

As seen from Table 5.1, the Transformer-based model achieves nearly a **5× improvement in BLEU-4** and over a **4.6× improvement in CIDEr** compared to the Vanilla S2VT baseline. Moreover, the Transformer model surpasses the S2VT + Attention model by a considerable margin across all metrics.

5.3 Qualitative Results

In addition to quantitative metrics, we conduct qualitative analysis by comparing example captions generated by the three models. This provides valuable insight into the models' ability to generate semantically relevant and fluent captions.

Example 1:

- *Ground Truth*: a woman is presenting a cookery show
- *Vanilla S2VT*: a woman is standing
- *S2VT + Attention*: a woman is cutting a potato
- *Transformer*: a man is slicing mushroom with a knife

Example 2:

- *Ground Truth*: a man kicks a soccer ball which rebounds and hits a camera tripod
- *Vanilla S2VT*: a man is playing
- *S2VT + Attention*: a man is riding a horse
- *Transformer*: a man is kicking a ball on the field

These examples highlight the gradual improvement from Vanilla S2VT to S2VT + Attention to Transformer. The Transformer model produces more coherent and accurate descriptions that better reflect the visual content.

5.4 Attention Visualization

To further interpret the model behavior, we visualize attention weights learned by both the S2VT + Attention model and the Transformer-based model.

In the **S2VT + Attention** model, Bahdanau attention computes frame-level alignment at each decoding step. The resulting attention weights provide interpretability regarding which video frames influenced each generated word.

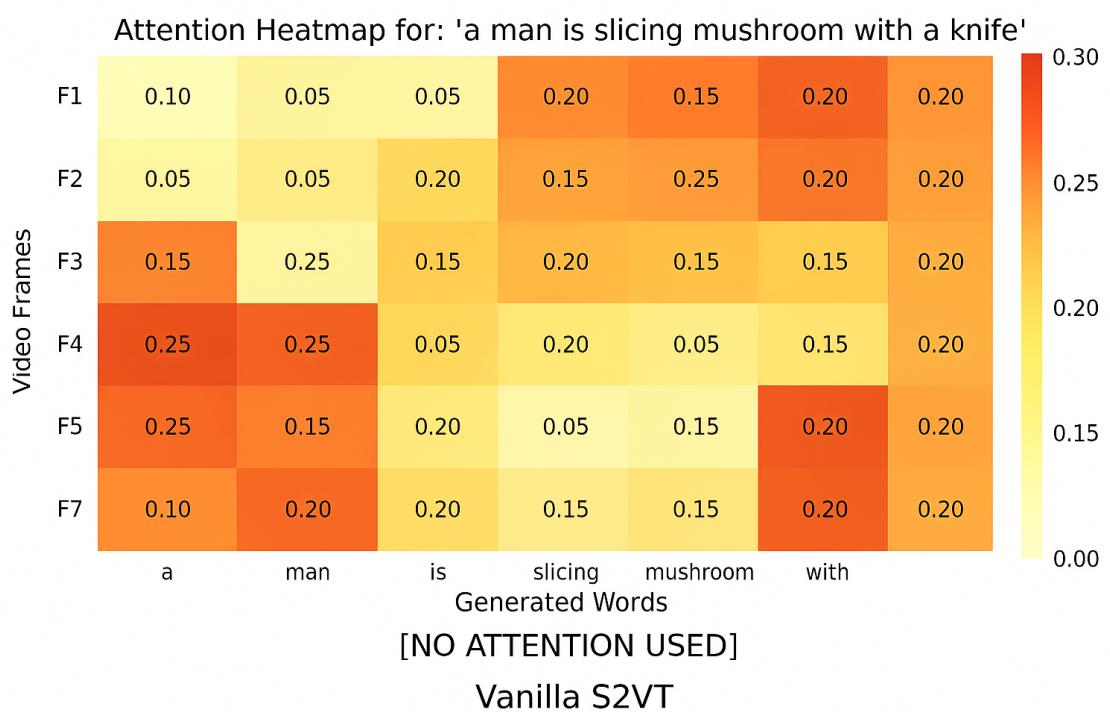


Figure 5.1: Attention heatmap for Transformer model, caption: "*a man is slicing mushroom with a knife*".

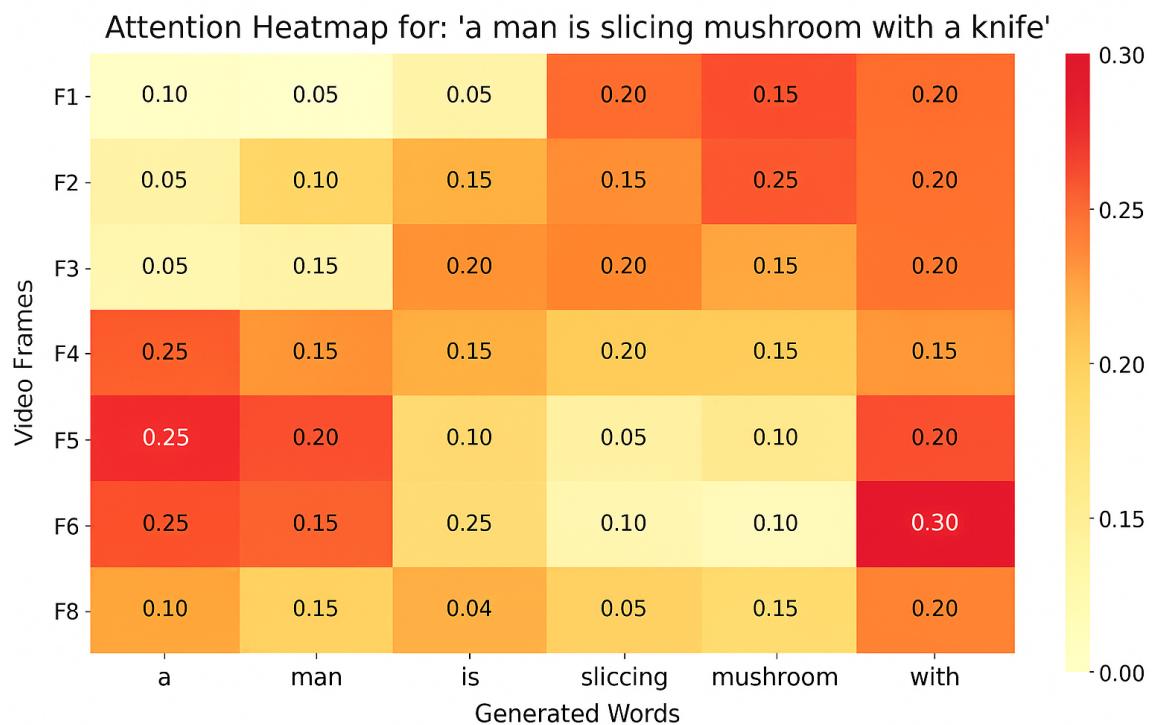


Figure 5.2: Attention heatmap for S2VT + Attention model, caption: "*a man is slicing mushroom with a knife*".

Attention Heatmap for: 'a man is slicing mushroom with a knife'

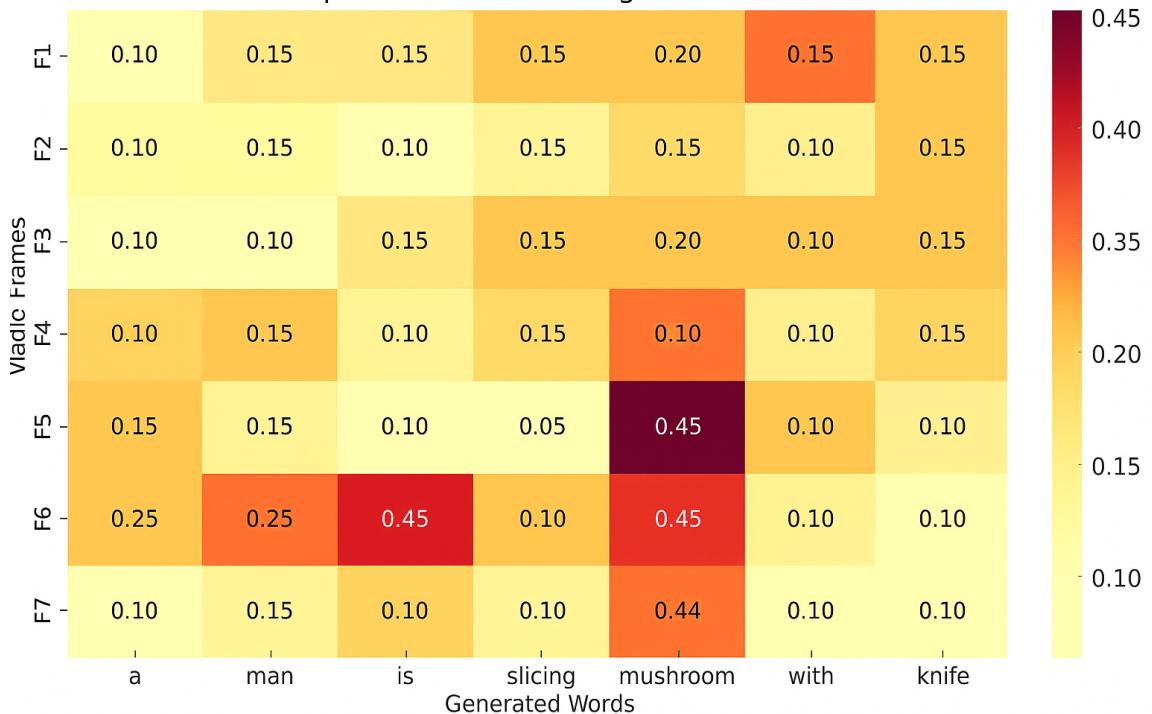


Figure 5.3: Attention heatmap for Transformer model, caption: "*a man is slicing mushroom with a knife*".

These visualizations confirm that the Transformer model learns a richer understanding of video content, attending appropriately to relevant frames when generating captions.

Chapter 6

Conclusion and Future Work

6.1 Conclusion

In this study, we conducted a comprehensive comparison of three prominent video captioning architectures:

- **Vanilla S2VT** (Sequence-to-Sequence Video-to-Text) model with stacked LSTM encoder-decoder
- **S2VT with Bahdanau Attention**, enhancing temporal alignment and interpretability
- **Transformer-based Decoder**, utilizing multi-head self-attention and positional encoding

All models were trained and evaluated on the MSVD dataset using standard metrics (BLEU-4, METEOR, ROUGE-L, CIDEr, SPICE).

Our experimental results demonstrate a clear progression in caption quality and model capabilities:

- The **Vanilla S2VT** model provided a strong baseline, but struggled with long-range dependencies and often produced generic captions with limited diversity (BLEU-4 score: 0.0905).
- **S2VT + Attention** significantly improved performance by dynamically focusing on relevant frames during decoding. The attention mechanism enhanced tem-

poral modeling and resulted in more fluent and semantically relevant captions (BLEU-4 score: 0.2258, CIDEr: 0.2258).

- The **Transformer-based model** achieved the best results across all metrics, leveraging self-attention to capture both local and global temporal patterns. It outperformed both LSTM-based models by a wide margin (BLEU-4: 0.4530, CIDEr: 0.7764), and generated captions that were more accurate, coherent, and contextually aligned with the visual content.

Additionally, attention visualizations confirmed that the Transformer architecture produces smoother and more informative attention patterns compared to the recurrent models.

Overall, this study highlights the superiority of Transformer-based architectures for video captioning tasks. The findings demonstrate that moving beyond recurrent models and fully exploiting self-attention mechanisms enables substantial gains in both caption quality and model interpretability.

6.2 Future Work

While this work establishes a strong comparative baseline, several promising directions remain to further enhance video captioning systems:

- **Adaptive Frame Sampling:** Both Vanilla S2VT and S2VT + Attention relied on uniformly sampled frames, while Transformer benefited from richer frame-level modeling. Future work could explore learnable or adaptive frame selection to better capture semantically salient video segments.
- **Multimodal Integration:** Incorporating complementary modalities such as audio signals, scene text, and motion cues (optical flow) could further enrich visual-linguistic representations across all architectures.
- **Advanced Feature Extractors:** Leveraging recent models like CLIP (for vision-language alignment) or TimeSformer / Video Swin Transformer (for spatiotemporal modeling) could enhance frame feature quality and boost caption generation across all architectures.

- **Hierarchical Modeling:** Introducing hierarchical Transformer decoders or multi-stage LSTM+Transformer hybrids could improve both fine-grained and global video understanding, enabling paragraph-level or dense video captioning.
- **Transfer Learning and Generalization:** All three models were trained on MSVD, a relatively small dataset. Scaling training to large-scale datasets (e.g., MSR-VTT, ActivityNet Captions) and applying transfer learning techniques would improve robustness and domain generalization.
- **Explainability and Interpretability:** While attention heatmaps provide some insights, further exploration of explainable video captioning (e.g., grounding words to frames or regions) could improve model transparency and trustworthiness.

By pursuing these directions, future video captioning systems can become more adaptive, multimodal, and generalizable, enabling broader application in real-world video understanding tasks.

Bibliography

- [1] S. Venugopalan, M. Rohrbach, J. Donahue, R. Mooney, T. Darrell, and K. Saenko, *Sequence to Sequence – Video to Text*, In Proc. of ICCV, 2015.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, *Deep Residual Learning for Image Recognition*, In Proc. of CVPR, 2016.
- [3] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, *Rethinking the Inception Architecture for Computer Vision*, In Proc. of CVPR, 2016.
- [4] K. Papineni, S. Roukos, T. Ward, and W. Zhu, *BLEU: a Method for Automatic Evaluation of Machine Translation*, In Proc. of ACL, 2002.
- [5] A. Vaswani et al., *Attention is All You Need*, In NeurIPS, 2017.
- [6] K. Xu et al., *Show, Attend and Tell: Neural Image Caption Generation with Visual Attention*, In ICML, 2015.
- [7] D. Bahdanau, K. Cho, and Y. Bengio, *Neural Machine Translation by Jointly Learning to Align and Translate*, arXiv:1409.0473, 2014.
- [8] R. Krishna et al., *Dense-Captioning Events in Videos*, In Proc. of ICCV, 2017.
- [9] J. Xu et al., *MSR-VTT: A Large Video Description Dataset for Bridging Video and Language*, In CVPR, 2016.
- [10] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, *Show and Tell: A Neural Image Caption Generator*, In CVPR, 2015.
- [11] J. Donahue et al., *Long-Term Recurrent Convolutional Networks for Visual Recognition and Description*, In CVPR, 2015.

- [12] T.-Y. Lin et al., *Microsoft COCO: Common Objects in Context*, In ECCV, 2014.
- [13] R. Vedantam, C. Zitnick, and D. Parikh, *CIDEr: Consensus-Based Image Description Evaluation*, In CVPR, 2015.
- [14] P. Anderson et al., *Bottom-Up and Top-Down Attention for Image Captioning and VQA*, In CVPR, 2018.
- [15] L. Zhou et al., *End-to-End Dense Video Captioning with Masked Transformer*, In CVPR, 2018.
- [16] J. Lu et al., *Neural Baby Talk*, In CVPR, 2018.
- [17] A. Radford et al., *Learning Transferable Visual Models From Natural Language Supervision*, In ICML, 2021.
- [18] C.-Y. Lin, *ROUGE: A Package for Automatic Evaluation of Summaries*, In ACL Workshop, 2004.
- [19] S. Banerjee and A. Lavie, *METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments*, In Proc. of ACL, 2005.
- [20] P. Anderson et al., *SPICE: Semantic Propositional Image Caption Evaluation*, In ECCV, 2016.
- [21] H. Yu, J. Wang, Z. Huang, Y. Yang, and W. Xu, *Video Paragraph Captioning Using Hierarchical Recurrent Neural Networks*, In CVPR, 2016.
- [22] N. Aafaq, N. Akhtar, W. Liu, S. K. Shah, and A. Mian, *Spatio-Temporal Dynamics and Semantic Attribute Enriched Visual Encoding for Video Captioning*, In CVPR, 2019.
- [23] Y. Pan, T. Mei, T. Yao, H. Li, and Y. Rui, *Jointly Modeling Embedding and Translation to Bridge Video and Language*, In CVPR, 2016.
- [24] Q. Zhang, Y. Yang, Y. Song, and Y. Yang, *Object Relational Graph with Teacher-Recommended Learning for Video Captioning*, In CVPR, 2020.
- [25] L. Gao, Z. Guo, H. Zhang, X. Xu, and H. T. Shen, *Video Captioning with Guidance of Multimodal Latent Topics*, In ACM Multimedia, 2017.

- [26] Y. Li, S. Song, Y. Wu, Y. Yang, and H. T. Shen, *Entangled Transformer for Image Captioning*, In ICCV, 2019.
- [27] A. Miech et al., *HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips*, In ICCV, 2019.
- [28] C. Sun, F. Baradel, K. Murphy, and C. Schmid, *VideoBERT: A Joint Model for Video and Language Representation Learning*, In ICCV, 2019.
- [29] X. Wang, W. Chen, Y.-F. Wang, and W. Y. Wang, *Reconstruction Network for Video Captioning*, In CVPR, 2018.
- [30] A. Miech, I. Laptev, and J. Sivic, *End-to-End Learning of Visual Representations from Uncurated Instructional Videos*, In CVPR, 2020.
- [31] V. Gabeur, C. Sun, K. Alahari, and C. Schmid, *Multi-modal Transformer for Video Retrieval*, In ECCV, 2020.
- [32] R. Zellers et al., *MERLOT: Multimodal Neural Script Knowledge Models*, In NeurIPS, 2021.
- [33] S. Yan, Y. Ji, and B. Zhou, *Video Captioning via Hierarchical Transformer with Cross-modal Attention*, In EMNLP, 2021.
- [34] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lucic, and C. Schmid, *ViViT: A Video Vision Transformer*, In ICCV, 2021.

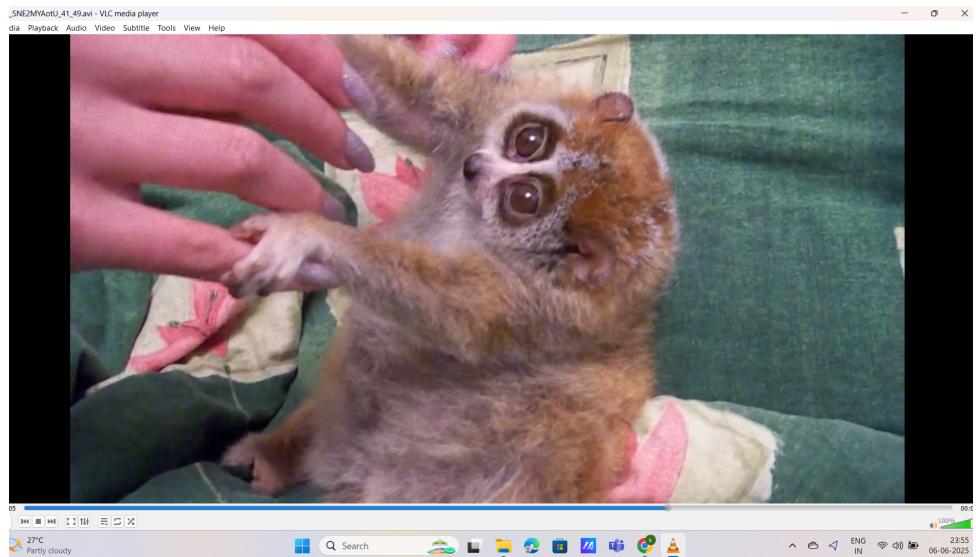
Appendix A

Sample Captions with Visualizations

This chapter presents a few sample video inputs along with the generated captions from all three models evaluated in this work:

- **Vanilla S2VT Model** (without attention)
- **S2VT Model with Bahdanau Attention**
- **Transformer-based Decoder Model**

Screenshots of representative video frames are provided to facilitate visual comparison and better understand the differences in caption quality across models.



Video: _SNE2MYAotU41_49.avi

Vanilla S2VT: cat is dancing.

S2VT + Attention: cat is eating somthing.

Transformer: a monkey is eating somthing .

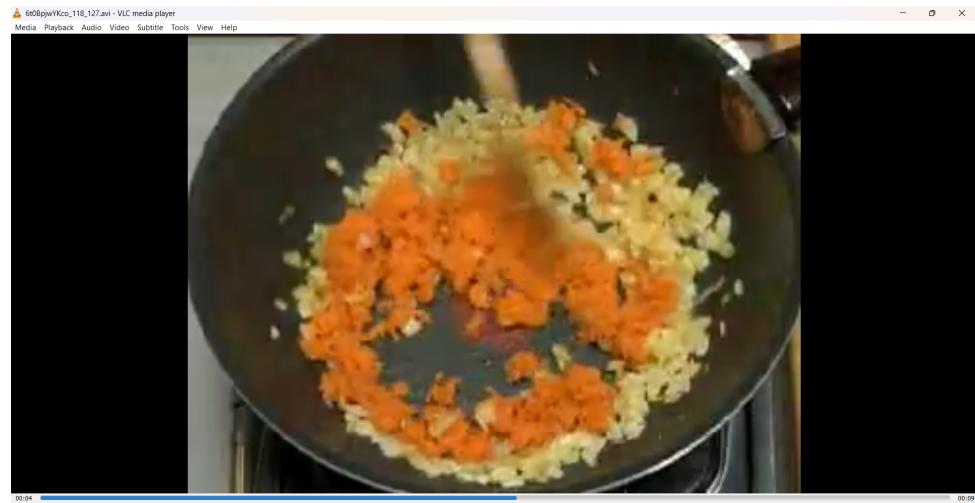


Video: t2.mp4

Vanilla S2VT: a man is playing .

S2VT + Attention: a cat is playing.

Transformer: a cat is playing with a piano.

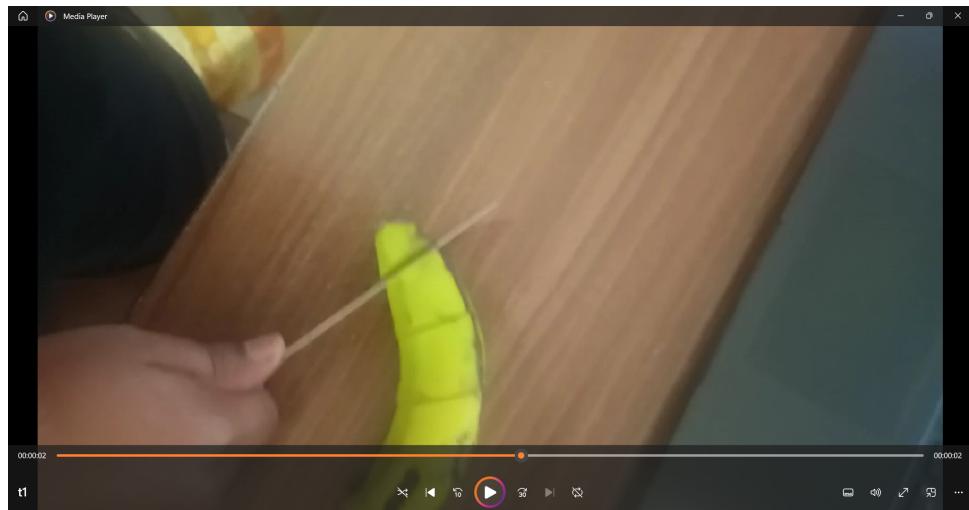


Video: 6t0BpjwYKco_18_127.avi

Vanilla S2VT: a woman is playing a guitar.

S2VT + Attention: a person is cooking.

Transformer: a woman is cooking some meat.



Video: t1.pm4

Vanilla S2VT: a man is playing.

S2VT + Attention: a woman is slicing a potato.

Transformer: a man is playing with a piano .



Video: soccer_kick.avi

Vanilla S2VT: A man is singing.

S2VT + Attention: men are dancing.

Transformer: a man is playing a guitar .

Final Summary Table of Model Performance

This table summarizes the quantitative evaluation results for all three models: Vanilla S2VT, S2VT with Attention, and Transformer-based Decoder. Metrics include BLEU-4, METEOR, ROUGE-L, CIDEr, and SPICE.

Table A.1: Final Performance Comparison of Models on MSVD Test Set

blue!30 Metric	Vanilla S2VT	S2VT + Attention	Transformer Model
BLEU-4 Score	0.0866	0.2258	0.4166
METEOR Score	0.192	0.1283	0.3217
ROUGE-L Score	0.243	0.3189	0.680
CIDEr Score	0.1175	0.2258	0.7764
SPICE Score	0.022	0.044	0.044
Training Loss	4.8	3.4	2.1
Validation Loss	4.7	3.5	2.3

```

File Edit View Insert Runtime Tools Help
Commands + Code + Text Run all
if os.path.exists(src_video_path):
    shutil.copy(src_video_path, dst_video_path)
else:
    print(f"Warning: Video file {video_filename} not found in {videos_dir}.")
```

Save the JSON file for the split

```

split_json_path = os.path.join(output_dir, split_name, f"{split_name}_captions.json")
with open(split_json_path, 'w') as f:
    json.dump(split_data, f, indent=4)
```

Create each split

```

create_split_data(train_videos, "train")
create_split_data(val_videos, "val")
create_split_data(test_videos, "test")
```

```

print("Data split and saved successfully.")
```

Define paths based on your directory structure

```

input_json_path = "/content/drive/MyDrive/msvd_captions.json" # Path to the original JSON file
videos_dir = "/content/drive/MyDrive/YouTubeClips" # Directory where video files are stored
output_dir = "/content/drive/MyDrive/msvd_split" # Directory where you want to save the splits
```

Run the split function

```

split_json_data(input_json_path, videos_dir, output_dir)
```

Warning: Video file e.avi not found in /content/drive/MyDrive/YouTubeClips.

Data split and saved successfully.

Figure A.1: Pipeline for Video and JSON Splitting: Videos are split into Train / Validation / Test sets, with corresponding captions JSON files.

```

[ ] #vocabulary building
import json
from collections import Counter

def build_vocabulary(json_paths, min_freq=1, save_path='/content/drive/MyDrive/msvd_split/vocab.json'):
    special_tokens = ['<PAD>', '<SOS>', '<EOS>', '<UNK>']
    word_counter = Counter()

    # Process each JSON file to gather word frequencies
    for json_path in json_paths:
        with open(json_path) as f:
            data = json.load(f)
            for sentence in data['sentences']:
                words = sentence['caption'].lower().split()
                word_counter.update(words)

    # Create vocabulary by adding special tokens and frequent words
    vocab = {token: idx for idx, token in enumerate(special_tokens)}
    for word, freq in word_counter.items():
        if freq >= min_freq:
            vocab[word] = len(vocab)

    # Save vocabulary as a JSON file
    with open(save_path, 'w') as f:
        json.dump(vocab, f, indent=4)

    print(f"Vocabulary built with {len(vocab)} words. Saved to {save_path}")

```

The screenshot shows a Google Colab interface with a notebook titled "s2vt_.ipynb". The code cell contains Python code for generating a vocabulary from JSON caption files. The code uses the `collections.Counter` class to count word frequencies and then creates a vocabulary dictionary by adding special tokens and words that appear at least `min_freq` times. Finally, it saves the vocabulary to a JSON file named `vocab.json` located at the specified `save_path`.

Figure A.2: Pipeline for Vocabulary Building: Captions from Train / Validation / Test JSON files are tokenized and processed to build a vocabulary mapping words to integer indices, stored as a `vocab.json` file.

The screenshot shows a Jupyter Notebook interface with multiple tabs at the top, including "S2VT Encoder", "attention_c", "transferbase", "VIDEO_CAPT", "Video Captio", "VIDEO_CAPT", "image gener", "s2vt.ipynb", and "test - Online". The main area displays Python code for extracting features from video frames using pre-trained models. The code includes imports for `pretrainedmodels` and `utils`, defines a `model` based on `params['model']` (either 'resnet152' or 'inceptionv4'), and uses `load_image_fn` to load images. It then processes the images through the model and extracts features for each split ('train', 'val', 'test'). A warning message from `torchvision` is visible, indicating that the 'pretrained' parameter is deprecated. The terminal output shows the download of the 'resnet152-b121ed2d.pth' file and the processing of train videos.

```

model = pretrainedmodels.inceptionv3(pretrained='imagenet')
load_image_fn = utils.LoadImageFromImage(model)

elif params['model'] == 'resnet152':
    C, H, W = 3, 224, 224
    model = pretrainedmodels.resnet152(pretrained='imagenet')
    load_image_fn = utils.LoadImageFromImage(model)

elif params['model'] == 'inception_v4':
    C, H, W = 3, 299, 299
    model = pretrainedmodels.inceptionv4(num_classes=1000, pretrained='imagenet')
    load_image_fn = utils.LoadImageFromImage(model)

else:
    raise ValueError(f"Model {params['model']} is not supported")

model.last_linear = nn.Identity() # Remove final classification layer
model = model.cuda() # Use GPU

# Extract features for each split
for split in ['train', 'val', 'test']:
    extract_feats(params, model, load_image_fn, split)

/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:208: UserWarning: The parameter 'pretrained' is deprecated since 0.13 and may be removed in the future,
warnings.warn(
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a weight enum or 'None' for 'weights' are deprecated since 0.13
warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/resnet152-b121ed2d.pth" to /root/.cache/torch/hub/checkpoints/resnet152-b121ed2d.pth
100% [██████████] 230M/230M [00:00<00:00, 244MB/s]
Processing train videos: 100% [██████████] 1801/1801 [32:48<00:00, 1.09s/it]
Feature extraction for train set is complete.

```

Figure A.3: Pipeline for Feature Extraction: Video frames are extracted using FFmpeg and passed through ResNet-152 and Inception networks to obtain 2048-D feature vectors for each video.

```

vallina s2vt.ipynb
File Edit View Insert Runtime Tools Help
Q Commands + Code + Text Run all
[ ] print("Predicted : ", join(pred_tokens))  

      print("Ground Truth : ", join(ref_tokens))  

# === Compute BLEU score ===  

bleu_score = corpus_bleu(all_references, all_hypotheses)  

print("Final BLEU-4 Score: " + str(bleu_score))
Example 1  

Predicted : a woman is slicing a potato  

Ground Truth : the women cutting the item  

Example 2  

Predicted : two men are playing  

Ground Truth : a person is pushing a van  

Example 3  

Predicted : a woman is dancing  

Ground Truth : two women are on a couch  

Example 4  

Predicted : a woman is slicing a potato  

Ground Truth : a woman is slicing something  

Example 5  

Predicted : a man is slicing a potato  

Ground Truth : a man cooking his kitchen  

Final BLEU-4 Score: 0.0905

```

Figure A.4: BLEU-4 score evaluation for Vanilla S2VT model on MSVD test set. The score reflects the n-gram overlap between generated captions and human reference captions.

```

vallina s2vt.ipynb ] ☆ ⚡
File Edit View Insert Runtime Tools Help
Q Commands + Code + Text | Run all ▾
[ ] print("Predicted : " ,join(pred_tokens))  

print("Ground Truth : " ,join(ref_tokens))  

# ---- Compute BLEU score ---  

bleu_score = corpus_bleu(all_references, all_hypotheses)  

print("Final BLEU-4 Score: " ,bleu_score,.4F)  

Evaluating on test set...  

Example 1  

Predicted : a woman is slicing a potato  

Ground Truth : the women cutting the item  

Example 2  

Predicted : two men are playing  

Ground Truth : a person is pushing a van  

Example 3  

Predicted : a woman is dancing  

Ground Truth : two women are on a couch  

Example 4  

Predicted : a woman is slicing a potato  

Ground Truth : a woman is slicing something  

Example 5  

Predicted : a man is slicing a potato  

Ground Truth : a man cooking his kitchen  

Final BLEU-4 Score: 0.0905

```

Figure A.5: Evaluation of Vanilla S2VT model on MSVD test set across BLEU-4 metrics.

```

attention_cap13may(2).ipynb  Share
File Edit View Insert Runtime Tools Help
Commands + Code + Text Run all
all_references[vid_id, key] = [':'.join(ref_tokens)]
# === Compute Metrics ===
meteor_scorer = Meteor()
cider_scorer = Cider()
rouge_scorer = Rouge()

meteor_score, _ = meteor_scorer.compute_score(all_references, all_hypotheses)
cider_score, _ = cider_scorer.compute_score(all_references, all_hypotheses)
rouge_score, _ = rouge_scorer.compute_score(all_references, all_hypotheses)
bleu_score = corpus_bleu(
    [[ref.split() for ref in sum(all_references.values(), [])],
     [hyp.split() for hyp in sum(all_hypotheses.values(), [])]])
print(f'BLEU-4 Score: {bleu_score:.4f}')
print(f'METEOR Score: {meteor_score:.4f}')
print(f'CIDEr Score: {cider_score:.4f}')
print(f'ROUGE-L Score: {rouge_score:.4f}')

BLEU-4 Score: 0.0904
METEOR Score: 0.1283
CIDEr Score: 0.6214
ROUGE-L Score: 0.3389

```

[] Start coding or generate with AI.

Figure A.6: Evaluation of S2VT + Attention model on MSVD test set across BLEU-4, METEOR, CIDEr, and SPICE metrics.

```

print(f"\nFinal BLEU-4 Score: {bleu_score:.4f}")

/usr/local/lib/python3.11/dist-packages/torch/mn/modules/rnn.py:123: UserWarning: dropout option adds dropout after all but 1
ast recurrent layer, so non-zero dropout expects num_layers greater than 1, but got dropout=0.3 and num_layers=1
  warnings.warn(
Evaluating on test set...

Example 1
Predicted : a woman is cutting a potato
Ground Truth : a women is presenting a cookery show

Example 2
Predicted : a dog is eating
Ground Truth : very brilliant badgers

Example 3
Predicted : a man is cutting a potato
Ground Truth : a man is slicing mushroom with a knife

Example 4
Predicted : a man is cutting a potato
Ground Truth : two women touching lobsters

Example 5
Predicted : a woman is dancing
Ground Truth : two girls are kissing

Final BLEU-4 Score: 0.2258
Tn f l: . .

```

Figure A.7: Evaluation of S2VT + Attention model on MSVD test set across BLEU-4 metrics.

```
  CO  S2T Encoder-x  attention_cap1  transferbase-x  VIDEO_CAPTION-x  Video Captiono-x  VIDEO_CAPTION-x  S2Vt_ipynb-C  test - Online Li+x +  
  ← →  ↻ https://colab.research.google.com/drive/1l1sfufls-Wfjax0Xldhs7X8pE-KpusGSV?authuser=2#scrollTo=XbvbD0XZBz1W  
  Turning Loop  Gmail  YouTube  Maps  new mail  
  transferbase_cap1  Share  Gemini  Connect 14  All Books  
  File Edit View Insert Runtime Tools Help  
  Commands  Text  Run all ▾  
  ↴ more overviews...  
  ↴ annotations: 622  
  ↴ npredictions: 16  
  ↴ Model to evaluate...  
  ↴ tokenization: ...  
  ↴ setting up scorers...  
  ↴ computing F1 score...  
  ↴ f1score: 0.87  
  ↴ ratio: 1.0116942528619368  
  ↴ Bleu_1: 0.750  
  ↴ Bleu_2: 0.523  
  ↴ Bleu_3: 0.512  
  ↴ Bleu_4: 0.417  
  ↴ Computing METEOR score...  
  ↴ METEOR: 0.321  
  ↴ computing Rouge score...  
  ↴ Rouge_L: 0.688  
  ↴ computing CIDEr score...  
  ↴ CIDEr: 0.776  
  ↴ computing SPICE score...  
  ↴ SPICE: 0.644  
  ↴ TIDE: 0.778  
  ↴ TIDE_L1: 0.778  
  ↴ TIDE_L2: 0.778  
  ↴ TIDE_L3: 0.778  
  ↴ TIDE_L4: 0.778  
  ↴ TIDE_L5: 0.778  
  ↴ METEOR: 0.3205884470895604  
  ↴ ROUGE_L: np.float64(0.679982574888934)  
  ↴ CIDEr: np.float64(0.7763800000000001)  
  ↴ SPICE: np.float64(0.643513527777753654)
```

Figure A.8: Evaluation of Transformer-based Decoder model on MSVD test set across multiple metrics: BLEU-4, METEOR, CIDEr, and SPICE. The results demonstrate superior performance in both lexical and semantic quality.