



VEERMATA JIJABAI TECHNOLOGICAL INSTITUTE

BTECH PROJECT

DEPARTMENT OF COMPUTER ENGINEERING & INFORMATION TECHNOLOGY

Dense Video Captioning

Ganadhish Acharekar
Akshat Shah
Arnav Shah
Saharsh Jain

Project Supervisor
Prof. Sandip T. Shingade

Abstract

The task of Dense Video Captioning (DVC), featured in the annual ActivityNet competition, aims to localize the different events in a given video and generate natural language descriptions for each of them. Most literature for this task involves utilizing CNN extracted video features and multiple stages of training. This project aims to utilize multiple modalities, namely audio and video, as opposed to previous single modality approaches, and we propose a model architecture that can be trained end-to-end. We utilize transformer architectures and efficient attention mechanisms for feature extraction, proposal generation as well as for captioning. We report xxxxx performance.

Contents

1	Introduction	7
1.1	Problem Background	8
1.2	Motivation	9
1.3	Problem Formulation	9
1.4	Objectives	9
2	Literature Review	11
2.1	Temporal Action Proposals	11
2.2	Transformers in Computer Vision	11
2.3	Dense Video Captioning	12
2.4	Identified Themes	15
2.4.1	Based on Training schemes	15
2.4.2	Based on Modalities	15
2.4.3	Based on Task ordering	15
2.5	Datasets	15
2.6	Evaluation Metrics	18
2.7	Research Literature Gaps	19
2.7.1	Less Explored	19
2.7.2	Use of Transformers	19
2.7.3	Multimodal Architecture	20
2.7.4	Efficient Attention Mechanisms	20
2.7.5	End-to-End Training	20
3	Methodology	21
3.0.1	Multimodal Feature Extraction	21
	Video Features	22
	Audio Features	22
3.0.2	Training	23
	Knowledge Distillation	23
4	Feature Extraction	25
4.1	Feature Extraction Methodology	25
4.1.1	Video Features Extraction Using ViViT	25
4.1.2	Audio Features Extraction using AST	27
4.2	Temporal Sensitivity of Features	27
4.2.1	Temporally Sensitive Pretraining of Proposed Feature Extractors	28
	Architecture	28
	Data Flow	28
	Loss functions	29
A	Supplementary Details	31

B Paper Summaries	33
B.1 Dense-Captioning Events in Videos	33
B.2 A Better Use of Audio-Visual Cues: Dense Video Captioning with Bi-modal Transformer . . .	35
B.3 Vision Transformer (ViT) - An image is worth 16x16 words: Transformers for image recognition at scale	38
B.4 ViViT: A Video Vision Transformer	41
B.5 AST: Audio Spectrogram Transformer	44
B.6 Temporally Sensitive Pretraining of Video Encoders for Localization Tasks	46
B.7 End-to-end Dense Video Captioning with Parallel Decoding	48
B.8 End-to-End Object Detection with Transformers (DETR)	51
B.9 Deformable DETR: Deformable Transformers for End-to-End Object Detection	53
B.10 Sparse DETR: Efficient End-to-End Object Detection with Learnable Sparsity	55

Chapter 1

Introduction

Machine learning is a subfield of *Artificial Intelligence*, which involves algorithms that enable automatic learning of tasks based on data. Unlike classical programming, which involves giving instructions and input data to a computer and producing output data, machine learning aims to produce the instructions or rules used for a task as the result, when the input and corresponding output data is provided to a computer.

Deep Learning, a subfield of Machine Learning has been growing popular in recent years, owing to their performance over traditional machine learning techniques without the need for intensive feature engineering by practitioners; deep learning models involve building up successive layers of increasingly meaningful representations automatically, allowing them to be rich and expressive with information. Different forms of deep learning models, including dense or fully connected layers, convolutional networks, recurrent networks and others have been applied to a variety of tasks which traditional algorithms have failed to address.

Convolutional networks (CNNs) have dominated the field of computer vision for quite some time. Convolutional networks consist of convolutional layers, pooling layers and fully connected layers. Out of these, convolutional layers are the core building blocks of CNNs, which apply a series of different *image filters* or *convolutional kernels* to an input image in a convolution operation and the outputs of different filters at a layer are stacked together. The filter weights are learned, so as to learn to extract features while making use of the spatial structure inherent in images for the given task.

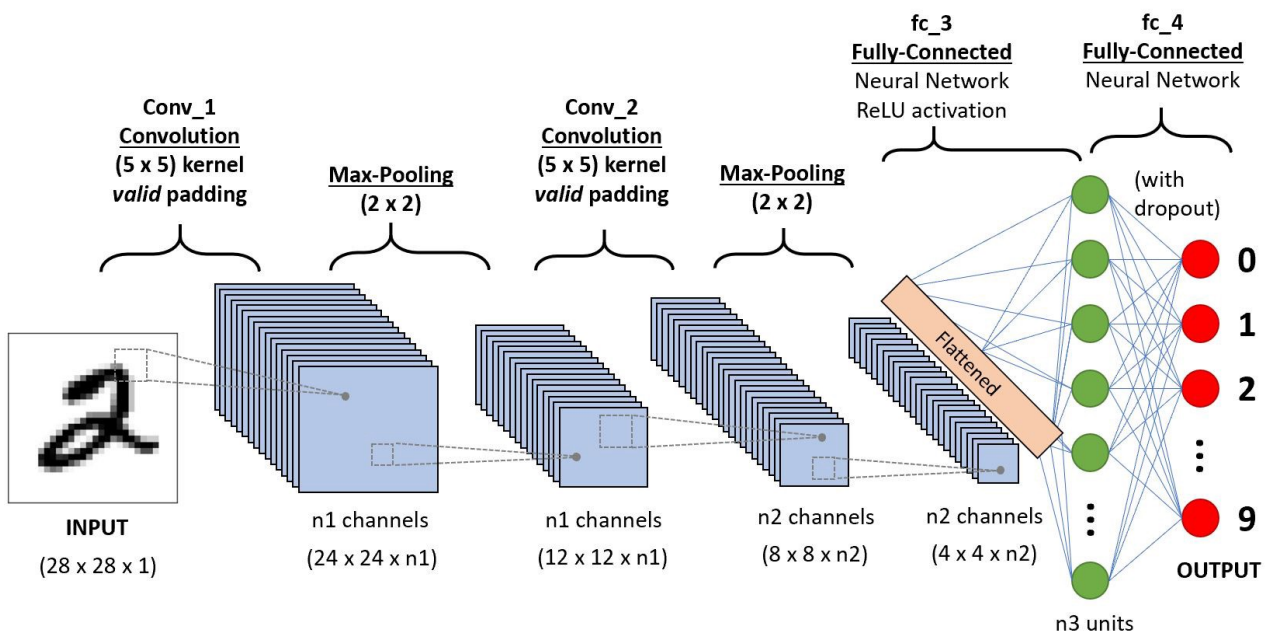


Figure 1.1: A Convolutional Network

Transformers were introduced by Vaswani *et al* in [1]. The transformer architecture is shown in Figure 1.2. It consists of encoder layers which use *multi-head self-attention* on input tokens followed by fully connected layers. The decoder layers use the output tokens and representations from encoder layers in a *multi-head cross-attention* block followed by fully connected layers. This architecture has been adapted by various other models to serve a variety of tasks. Transformers are gaining traction in a variety of tasks, including sequence to sequence tasks like machine translation, natural language generation and understanding, and recently, computer vision.

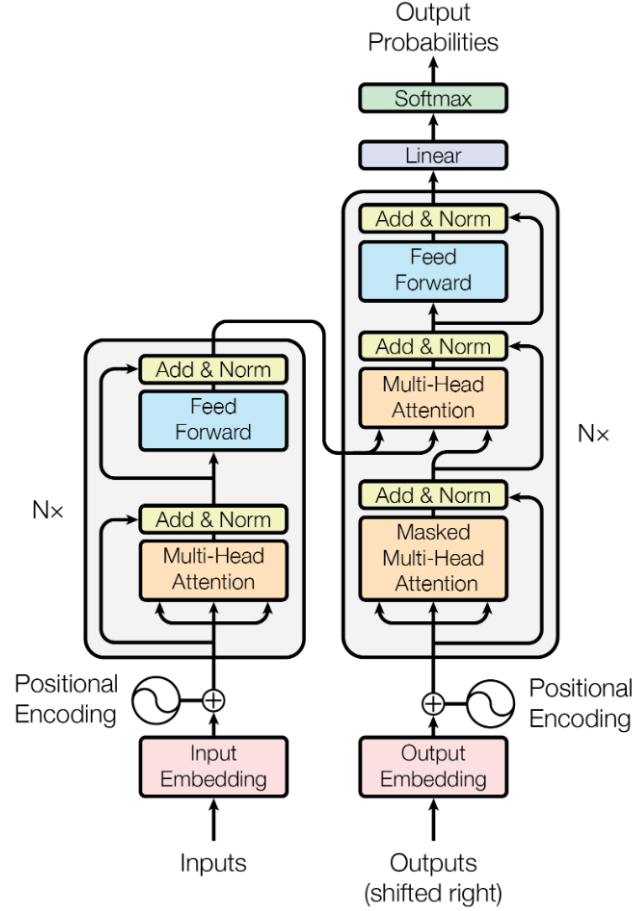


Figure 1.2: The Transformer architecture. Image courtesy [1]

Key Terms: Dense Video Captioning, Transformers, Temporally Sensitive Pretraining, Sparse Deformable Attention, Context

1.1 Problem Background

There has been a significant progress in the fields of *Computer Vision* (CV) and *Natural Language Processing* (NLP) and other artificial intelligence problems in recent decades. In computer vision, tasks like object detection and classification are well-explored, while in natural language processing, sequence to sequence tasks like machine translation have been well worked upon. The most promising solutions to these tasks come from deep learning methodologies, outperforming techniques like rule-based systems and other traditional machine learning algorithms. Owing to the rich feature representation capabilities of different kinds of neural networks, tasks that are more natural to humans are becoming solvable by machines, especially perceptive and understanding tasks.

The task of image captioning combines the fields of CV and NLP. Video action recognition, which involves classifying the action in a short video clip, is a task that can be considered as the representative task for video understanding. Taking one step further from action recognition and image captioning, the task of video

captioning involves labelling a video clip with a single natural language sentence. However, for long videos, a single sentence is not enough to describe all events that may occur in a video. Thus, the task of **Dense Video Captioning (DVC)** was introduced [2] in 2017. It involves generating natural language descriptions for each of the multiple events that may occur in a video. Krishna *et al* proposed a new dataset, ActivityNet Captions, which has segment annotations for untrimmed videos. Each annotation defines a start time, end time and a natural language description for each event segment in the videos.

1.2 Motivation

Describing a short video in natural language is trivial for humans, but a challenging task for machines. Automatic video description involves understanding of many entities and the detection of their occurrences in a video employing computer vision techniques. These entities include background scene, humans, objects, human actions, human-object interactions, human-human interactions, other events, and the order in which events occur. All this information should then be articulated using comprehensible and grammatically correct text, employing NLP techniques[3]. Indeed, the task of dense video captioning can be thought of as the culmination of all the perceptive learning tasks of computer vision and the natural language generation task.

As we move towards a digital age, more and more content is generated in the form of multimedia, particularly videos. The applications of dense video captioning are in tasks such as:

- Video summarization and highlights generation, e.g. in fields such as sports and entertainment
- Video retrieval via search and indexing on its dense captions. This can help in a multitude of areas such as experiments, surveillance and security, sports, education, et cetera.
- Video segment localization queries can be better handled by using dense event captions generated on a bulk of videos. Indeed, this can open the door to automatic interpretation of videos, a kind of data that is treated as unstructured binary data today.
- Increased accessibility to the visually impaired can be achieved with the combined use of DVC and the much more mature task of text-to-speech (TTS).

1.3 Problem Formulation

The task of dense video captioning involves two sub-tasks (1) temporal localization of events in a video and (2) describing the localized events in natural language. Formally, the task is defined as follows:

Given an input video $\mathbf{v} = \{v_1, v_2, \dots, v_T\}$ where v_i represents i^{th} video frame in temporal order, the target of dense video captioning task is to output a set of sentences $\mathbf{S} = \{s_1, s_2, \dots, s_{N_S}\}$ where N_S is the number of sentences and $s_i = \{t_i^{start}, t_i^{end}, \{w_j\}\}$ consists of start and end timestamps for each sentence described with set of words from a vocabulary set $w_j \in \mathbf{V}$.

Most of the architectures for dense video captioning comprises two components (1) Proposal generation module and (2) Captioning module. The modules can be trained in different ways like separate training, alternative training or in a end-to-end manner. The task of the proposal generation module is to take as input video frames v and output event proposals $P = \{t_i^{start}, t_i^{end}, confidence_i\}$. Depending on the architecture, the proposal generation module can also output additional parameters like in [2, 4, 5] which can be utilized by captioning module. The proposals can also be in the form of offsets (center and length) as in [4, 6, 7, 8, 9]. The captioning module outputs descriptions for each of the proposed event in usually a single sentence.

1.4 Objectives

- **End to end training:** Develop an end-to-end framework for dense captioning of long untrimmed videos. The task of dense video captioning consists of two sub tasks, namely temporal proposal generation and proposal captioning. These two tasks are highly dependent on one another and training them independently or in a weakly dependent way like alternate training can lead to poor correlation between

proposals and captions. The generated proposals should be able to generate good captions and the loss from generated captions should be able to backpropagate through the proposal generation module for correcting the timestamps. This backpropagation is only possible through an end-to-end training mechanism with strong relation between the two modules.

- **Multimodal input:** Utilize both video frame (visual) and audio features. Multimodal learning suggests that when a number of inputs – visual, auditory, kinaesthetic – are being used during learning, the model gets better understanding of the input data. Most of the current methods only utilize the visual modality from the videos. Taking inspiration from human perception system, which utilizes much more modalities than just the visual one for understanding the surrounding and making decisions, the model should exploit both audio and video modalities to take decisions on the proposal timestamps and captions.
- **Contextual learning:** Use temporal context while generating captions. The captions of the video are highly dependent on the time interval they correspond to. An event depends on the past, present and future events. In other words, there is existence of both inter-event and intra-event dependencies. The model should utilize these dependencies while predicting caption for an event.

Chapter 2

Literature Review

2.1 Temporal Action Proposals

The task of Temporal Action Proposals involves generating temporal segments in long untrimmed video. Given a video, the model should be able to output temporal intervals (timestamps) that have a possibility of containing an action. Buch *et al* [10] used Single-Stream Temporal Action Proposal (SST) that does not require division of input video into a number of short clips or use of overlapping sliding windows. The method uses a GRU-based sequence encoder which processes the video in a single pass and at each time step outputs proposals and their confidences by utilizing the C3D encoded features of the video frames. The model however had an issue of predicting an output segment which contains several ground-truth segments. Lin *et al* [11] develop a Boundary-Matching Mechanism for efficiently evaluating the proposals from a bottom-up generators and give reliable confidence scores.

2.2 Transformers in Computer Vision

Transformers [1] have been widely used in applications requiring processing of sequential data like language or speech. Various solutions have been coming forward to utilize this generalized architectures in the domain of computer vision. Dosovitskiy *et al* [12] used only transformers (called Vision Transformers or ViT), without any CNN backbone for the task of image classification. The model divided an image into fixed-size patches and pass them along with an extra learnable class embedding token as input to a transformer encoder. The output of the class embedding token is then feed to a MLP head to get the output object class. The approach suffers from the issues of less scalability and larger running time, since the number of tokens is highly dependent on the input resolution of the image. Liu *et al* [13] tackles these issues by introducing a hierarchical transformer called Swin-Transformer. The method uses shifted windows to limit the attention computation to non-overlapping local windows. It decreases the quadratic complexity of ViT to linear computational complexity with respect to the input image size. Fayyaz *et al* [14] introduces a parameter free Adaptive Token Sampling (ATS) module which can be plugged into any ViT models and reduce the GFLOPs without any additional training. The proposed method decides and selects the right amount of tokens to attend to.

Carion *et al* [15] develop DETection TRansformer (DETR) for the task of object detection or panoptic segmentation in images. The method completely eradicates the need of developing various hand-designed features like anchors, non-max suppression thresholds, etc. It is an end-to-end architecture viewing object detection task as set prediction task. The model tokenizes the input image using a CNN backbone (ResNet-50) and passes it to the transformer encoder along with positional embeddings. The output of encoder is given to the decoder along with learnable positional embeddings referred to as object queries. The decoder outputs the query features which are passed to a shared Feed Forward Network (FFN) that predicts the object bounding box and class. They use Hungarian Bipartite Matcher to compute the loss between ground truth and prediction. DETR inspite of comparable performance with other state-of-the-art CNN architectures, faced issues of slower convergence because of global attention in encoder and inaccuracy in detection of smaller objects. A number of approaches were introduced for handling this issue. Zhu *et al* [16] proposed Deformable DETR which answers

both the issues by using Multi-scale Feature Maps of the input image and attending to limited number of tokens in encoder. The method also introduced variations for better accuracy, namely iterative bounding box refinement and two-stage Deformable DETR. The model required 10x faster convergence and is 1.6x faster than the original DETR. This however does increase the input tokens by 20x because of multi-scale features. Roh *et al* [17] make Deformable DETR further efficient by selectively updating only the tokens that are expected to be referenced by the DETR decoder. This is done by learning a scoring network for predicting binarized Decoder cross-Attention Map (DAM). The model also improves the gradient flow throughout the transformer by introducing encoder-auxiliary loss. Sparse-DETR along with Swin Transformer Backbone increases the FPS by 38% as compared to Deformable DETR.

2.3 Dense Video Captioning

The problem of dense video captioning was introduced by Krishna *et al* in [2] by proposing ActivityNet Captions dataset. Their architecture involved a proposal module and captioning module. The proposal module was inspired from DAPs [18], while the captioning module incorporated LSTM with contextual features along with event feature as inputs. The architecture was able to detect events and generate captions of the video in a single pass without the need of a time consuming sliding window approach. However, since the features were dependent on the end location of the event, the model generated same captions for events ending at same timestamp.

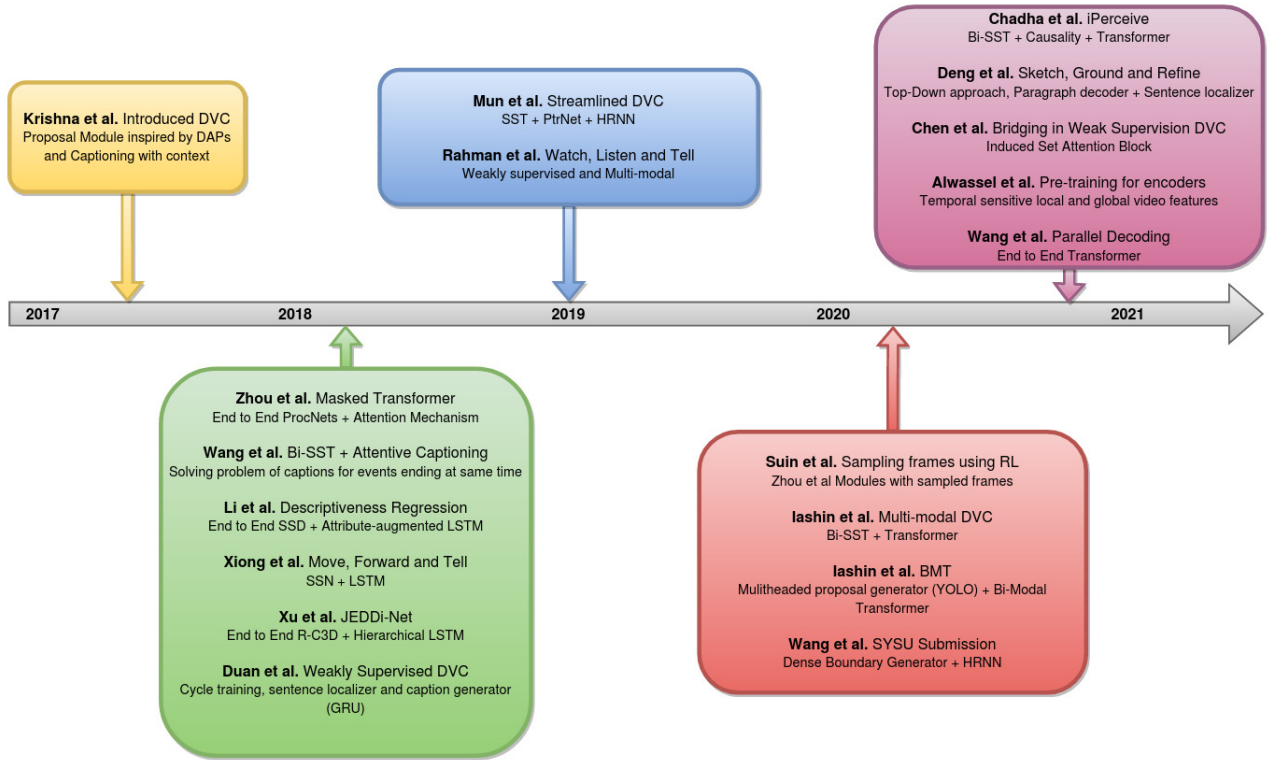


Figure 2.1: Evolution of dense video captioning methods over time.

Following the release of ActivityNet Captions dataset in 2017, many researchers were able to surpass the results of baseline model and achieve state-of-the-art. Zhou *et al* [9] addressed the problem of little influence of language descriptions on event proposals if the two modules are trained separately. They introduced an end-to-end masked transformer for propagating captioning error to proposal module for better performance. Furthermore, they proposed a self-attention mechanism for learning long-range dependencies in video. The proposal module was based on ProcNets [19]. The caption decoder employed a differentiable proposal mask to account for features in the respective event. Wang *et al* [5] employed Bi-SST as their proposal generator to

account for past and future context information. The captioning module consisted of attentive fusion of context features, the weights of which were decided by a context gating mechanism. The architecture selected final captions based on joint ranking method which accounted for both proposal and caption confidence. Li *et al* [4], inspired by object localization networks like [20, 21], presented an end-to-end model with descriptiveness regression. An attribute-augmented LSTM network optimized using reinforcement learning was used for captioning module. Xiong *et al* [22] strived to generate relevant, coherent and concise descriptions using SSN [23] for event localization and LSTM for event selection and caption generation. Reinforcement learning with sentence-level reward is used to train the captioning network. Xu *et al* [7] proposed JEDDi-Net, an end-to-end architecture incorporating visual and language contexts. Segment Proposal Network inspired from R-C3D [24] is used for proposal generation. Hierarchical LSTM with caption-level controller network and word-level sentence decoder is used for caption generation. The proposal features are represented using 3D Segment-of-Interest Pooling. Duan *et al* [25] introduced weak supervision training with no need of temporal annotations of events. They trained sentence localizer and caption generator in cyclic manner minimizing the reconstruction loss. However, the method struggles to detect beginning of events properly.

Mun *et al* [26] tackles the challenge of coherent captioning by considering temporal dependency between events. They used SST for event proposals, PtrNet for event sequence generation and HRNN for captioning. Rahman *et al* [27] utilized weak supervision method from [25] and were the first to try multi-modal approach for dense video captioning. They show how audio alone can be competitive to the previous visual based results. The paper also discusses various methods to encode audio (MFCC, CQT, SoundNet) and context fusion techniques for multiple modalities (Multiplicative Mixture, Multi-modal context fusion, MUTAN). The architecture suffered from proposal localization accuracy due to weak supervision. Furthermore, the results are also affected due to unavailability of part of the dataset as some videos are not available on their respective urls.

Suin *et al* [28] aimed to reduce computational cost by processing fewer frames. They used deep reinforcement-based approach to describe multiple events in a video by watching a portion of the frames. The event proposal and captioning modules were inspired from [9]. Iashin *et al* [29] shows the importance of audio and speech modalities alongside visual features for dense video captioning task. They employ a Bi-SST for proposal and the captioning module consists of a Transformer architecture with three blocks: an encoder, decoder and generator. The model was able to achieve better performance than the then existing methods despite of unavailability of full dataset for multiple modalities. Iashin *et al* [6] utilized audio and video with Bi-modal Transformer for captioning. The proposal generator consisted of multiheaded method, inspired from YOLO object detector [30]. Their ablation analysis depicted stronger contribution of visual cues alone than audio cues alone. However, both modalities combined gave better results. Wang *et al* [31] adapt DBG [32] for temporal event proposals alongwith ESGN [26] for candidate subset selection. The captioning module consists of an encoder-decoder architecture with CMG (cross-modal gating) block to adaptively balance the visual and linguistic information.

Chadha *et al* [33] proposed to handle cognitive error (causality between events) and incorrect attention (attending to objects in the frame). Their end-to-end model consisted of Bi-SST for proposals, Common-Sense Reasoning for causality learning and Transformer based architecture [29] for captioning. The common-sense reasoning module employed a borrow-put experiment for determining dependency between events and generate context-aware features. Deng *et al* [8] introduced a top-down approach, reversing the usual detect-then-describe method. The architecture first generates a multi-sentence paragraph to describe the whole video and then localize each sentence for events. The captions are then refined using dual-path cross attention module. The top-down approach increased coherency in captions. Chen *et al* [34] worked on closely bridging event localization and captioning modules for weakly supervised learning. The Induced Set Attention Block helps the captioner to learn highly abstracted global structure of the video. The method suffers from detecting visually small concepts/objects. Alwassel *et al* [35] introduce a supervised pre-training paradigm for temporal action localization. The paradigm also considers background clips and global video information, to improve temporal sensitivity. Wang *et al* [36] formulated dense video captioning as a set prediction task and introduced an end-to-end dense video captioning framework with parallel decoding (PDVC). PDVC adopts the vision transformer to learn attentive interaction of different frames. Two prediction heads run in parallel over query features, leveraging the mutual benefits between two tasks of localization and captioning.

Paper	Set	Ground Truth Proposals				Learnt Proposals			
		M	C	B@3	B@4	M	C	B@3	B@4
Krishna <i>et al</i> [2] DCE	Validation	8.88	25.12	4.09	1.60	5.69	12.43	1.90	0.71
	Test	9.46	24.56	7.12	3.98	4.82	17.29	3.86	2.20
Zhou <i>et al</i> [9] Masked Transformer	Validation	11.16	47.71	5.76	2.71	4.98	9.25	2.42	1.15
	Test	10.12				10.12			
Wang <i>et al</i> [5] Bi-SST	Validation	10.89				5.86	7.99	2.55	1.31
	Test					9.65			
Li <i>et al</i> [4] Descriptivness Regr.	Validation	10.33	26.26	4.55	1.71	6.93	13.21	2.27	0.74
	Test					12.96			
Xu <i>et al</i> [7] JEDDi-Net	Test			4.06	1.63	8.58	19.88	4.06	1.63
Mun <i>et al</i> [26] Streamlined DVC	Validation	13.07	43.48	4.41	1.28	8.82	30.68	2.94	0.93
	Test					8.19			
Duan <i>et al</i> [25] Weakly Supervised DCE	Test			2.62	1.27	6.3	18.77	2.62	1.27
Rahman <i>et al</i> [27] Watch, Listen, Tell	Validation *	7.23	25.36	3.04	1.46	4.93	13.79	1.85	0.9
Suin <i>et al</i> [28] Efficient framework for DVC	Validation			2.87	1.35	6.21	13.82	2.87	1.35
Iashin <i>et al</i> [29] Multi-modal DVC	Validation *	11.72		5.83	2.86	7.31		2.6	1.07
Iashin <i>et al</i> [6] BMT	Validation *	10.90		4.63	1.99	8.44		3.84	1.88
Xiong <i>et al</i> [22] Move Forward and Tell	Validation			2.84	1.24	7.08		2.84	1.24
Wang <i>et al</i> [31] SYSU	Validation	14.85				10.31			
	Test					9.28			
Chadha <i>et al</i> [33] iPerceive	Validation *	12.27		6.13	2.98	7.87		2.93	1.29
Deng <i>et al</i> [8] Sketch, Ground and Refine	Validation				1.67	9.37	22.12		1.67
Chen <i>et al</i> [34] Towards Bridging EC-SL	Validation			2.78	1.33	7.49	21.21	2.78	1.33
Alwassel <i>et al</i> [35] TSP with BMT	Validation			4.16	2.02	8.75		4.16	2.02
Wang <i>et al</i> [36] Parallel decoding	Validation *	11.26	53.65		3.12	8.08	28.59		1.96

Table 2.1: Performance comparison of previous methods on the ActivityNet Captions dataset (* some videos unavailable)

2.4 Identified Themes

Dense video captioning can be decomposed into two parts: event localization and event description. Existing research methodologies can be grouped into different categories based on training methods, modalities used and ordering of tasks.

2.4.1 Based on Training schemes

1. **Independent training:** Training the proposal and captioning module independently, with no loss propagation or interaction between the two.
2. **Alternate training:** alternate between i) training the proposal module only and ii) training the captioning module on the positive event proposals while fine-tuning the proposal module.
3. **End-to-End:** Training the proposal and captioning module together, with caption information being able to affect and improve localization capability in the video.
4. **Weakly Supervised:** Training the model without ground-truth annotations of temporal events. This method assumes one-to-one correspondence between i.e. each caption describes one temporal event, and each temporal event has one caption.

2.4.2 Based on Modalities

1. **Uni-modal:** Utilizes only one type of input feature. For example, only visual features.
2. **Multi-modal:** Utilizes more than one type of input features. For example, combinations of visual, audio and speech features.

2.4.3 Based on Task ordering

1. **Top-Down:** *localize-then-describe*. Predicts the temporal intervals and then caption them. May include optional postprocessing such as ranking proposals, extracting top-k proposals based on confidence scores, etc.
2. **Bottom-Up:** *describe-localize-refine*. Describes the entire video into a paragraph using Video Paragraph Generation models. Then localize each sentence from the paragraph over the video. The localized boundaries are then used to re-caption the event.

2.5 Datasets

2.5.1 Textually Annotated Cooking Scenes (TACoS)

TACoS is a subset of MP-II Composites. It was then processed further to provide coherent textual descriptions for high-quality videos. The TACoS dataset was created by filtering through MP-II Composites and limiting the activities to those that involve the manipulation of cooking ingredients and have at least four videos for the same activity. As a result, TACoS includes 26 fine-grained cooking activities spread across 127 videos. Twenty different textual descriptions were collected for each video. The dataset consists of 11,796 sentences with 17,334 action descriptions. The dataset also allows for the alignment of sentences describing activities by obtaining approximate time stamps for the beginning and end of each activity.

2.5.2 TACoS-MultiLevel

The TACoS corpus was used to collect TACoS-MultiLevel. Three levels of descriptions were collected for each video in the TACoS corpus: (1) a detailed description of the video with no more than 15 sentences per video; (2) a short description with 3-5 sentences per video and (3) a single sentence description of the video. The data is annotated using tuples such as object, activity, tool, source, and target, with a person always serving as the subject.

2.5.3 Microsoft Video Description (MSVD)

MSVD is made up of 1,970 YouTube clips with human annotated sentences. To avoid bias from lexical choices in the descriptions, the audio is muted in all clips. Furthermore, videos with subtitles or overlaid text were removed during the dataset formulation quality control process. Each video in this dataset is typically between 10 and 25 seconds long and focuses on one activity. The dataset contains multilingual (Chinese, English, German, and so on) human generated descriptions.

2.5.4 Montreal Video Annotation Dataset (M-VAD)

M-VAD is based on the Descriptive Video Service (DVS) and includes 48,986 video clips from 92 films. Each clip lasts an average of 6.2 seconds, and the total time for the entire dataset is 84.6 hours. There are a total of 55,904 sentences, with only a few clips associated with more than one sentence. The dataset's vocabulary contains approximately 17,609 words (Nouns-9,512, Verbs-2,571, Adjectives-3,560, Adverbs-857). The dataset is divided into 38,949, 4,888, and 5,149 video clips for training, validation, and testing.

2.5.5 MPII Movie Description Corpus (MPII-MD)

MPII-MD includes transcribed audio descriptions from 94 Hollywood films. These films are divided into 68,337 clips with an average length of 3.9 seconds and 68,375 sentences, for a total of nearly one sentence per clip. Each clip is paired with one sentence taken from the film's script and audio description data. The Audio Descriptions (ADs) were first gathered by retrieving the audio streams from the film using the online services MakeMkv 1 and Subtitle Edit 2. The transcribed texts were then time stamped and aligned with the associated spoken sentences. The total duration of the dataset videos is nearly 73.6 hours, and the vocabulary size is 653,467 words.

2.5.6 MSR Video-to-Text (MSR-VTT)

MSR-VTT has a large collection of open domain videos for video captioning tasks. It is made up of 7180 videos divided into 10,000 clips. The clips are divided into 20 distinct categories. The dataset is split into 6513 training videos, 497 validation videos, and 2990 test videos. Each video contains 20 reference captions that have been annotated by AMT employees. This is one of the largest video captioning datasets in terms of the number of clips with multiple associated sentences. This dataset contains audio information in addition to video content, which could be used for multimodal research.

2.5.7 ActivityNet Captions

It contains 100k dense natural language descriptions of approximately 20k videos from ActivityNet, totaling 849 hours. Each description contains 13.48 words on average and covers approximately 36 seconds of video. Every video has multiple descriptions, and when combined, these descriptions cover 94.6 percent of the content in the video. Furthermore, the dataset's 10 percent temporal overlap makes it particularly interesting and challenging for studying multiple events that occur at the same time.

2.5.8 ActivityNet Entities

It is the first video dataset with grounding and annotations for entities. This dataset is based on the ActivityNet Captions dataset's training and validation splits, but with different captions. In this dataset, noun phrases (NPs) from video descriptions have been grounded to video frame bounding boxes. The dataset includes 14281 annotated videos, 52k video segments with at least one noun phrase annotated per segment, and 158k annotation-enhanced bounding boxes. The dataset uses a training set of 10,000 captions, similar to ActivityNet Captions. The validation set of ActivityNet Captions, on the other hand, is randomly and evenly divided into ANet-Entities validation (2.5k) and testing (2.5k) sets.

2.5.9 YouCook

It is made up of 88 YouTube cooking videos of various people preparing various recipes. Most of the videos have a different background (kitchen / scene). The dataset is divided into six different cooking styles, such as grilling, baking, and so on. The training set for machine learning contains 49 videos, while the test set contains 39 videos. The dataset's object categories include "utensils," "bowls," and "food," among others.

2.5.10 YouCook2

The YouCook-II Dataset contains 2000 videos that are evenly distributed across 89 recipes. The cooking videos are sourced from YouTube and include all of the challenges of open domain videos, such as camera position changes, camera motion, and changing backgrounds. The entire dataset has a play time of 175.6 hours and a vocabulary of 2600 words. The videos are further broken down into 3-16 segments per video, with an average of 7.7 segments per video elaborating on procedural steps. The length of each segment ranges from 1 to 264 seconds. Each video lasts an average of 316 seconds and can last up to 600 seconds. The dataset is randomly divided into train, validation, and test sets in the following proportions: 66%:23%:10%.

2.5.11 MP-II Cooking

The Cooking dataset from the Max Planck Institute for Informatics (MP-II) consists of 65 fine-grained cooking activities performed by 12 participants while preparing 14 dishes such as fruit salad and cake. Among the 65 cooking activities are "wash hands," "put in bowl," "cut apart," "take out from drawer," and so on. A "background activity" is generated when the person is not in the scene for 30 frames (one second) or is performing an unannotated activity. The dataset contains 44 videos (888,775 frames) with an average length of approximately 600 seconds per clip. The dataset contains a total of 8 hours of video play time and 5,609 annotations.

2.5.12 VideoStory

VideoStory is a dataset of 20k social media videos with multi-sentence descriptions. This dataset is intended to address the generation of story narration or description for long videos that cannot be adequately illustrated with a single sentence. Each video is accompanied by at least one paragraph. 4.67 sentences are temporally localised on average per paragraph. The dataset contains 26245 paragraphs and 123k sentences, with an average of 13.32 words per sentence. Each paragraph, on average, covers 96.7 percent of the video content. The dataset has approximately 22% temporal overlap between co-occurring events. The dataset includes 17908 videos for training, 999 videos for validation, and 1011 videos for testing, as well as a blind test set of 1039 videos.

2.5.13 Charades

It has 9848 videos of everyday indoor household activities. Videos are shot in 15 different indoor scenes using only 46 objects and 157 action classes. The dataset includes 66500 annotations that describe 157 actions. It also assigns 41104 labels to each of its 46 object classes. It also includes 27847 descriptions for all of the videos. The videos in the dataset have an average duration of 30 seconds and depict daily life activities. For training and testing purposes, the dataset is divided into 7985 and 1863 videos, respectively.

2.5.14 Video Titles in the Wild (VTW)

It has 18100 video clips with an average duration of 1.5 minutes. Each clip is described in a single sentence. However, it incorporates a diverse vocabulary, with one word appearing in no more than two sentences on average across the entire dataset. In addition to the single sentence per video, the dataset includes accompanying descriptions (known as augmented sentences) that describe information not present in the clip's visual content. The dataset is intended for video title generation rather than video content description, but it can also be used for language-level comprehension tasks such as video question answering.

2.5.15 Kinetics

It is a collection of large-scale, high-quality datasets of URL links to up to 650,000 video clips that, depending on the dataset version, cover 400/600/700 human action classes. Human-object interactions, such as playing instruments, are included in the videos, as well as human-human interactions, such as shaking hands and hugging. There are at least 400/600/700 video clips in each action class. Each clip is approximately 10 seconds long and human annotated with a single action class.

Dataset	Domain	#classes	#videos	Avg len	#clips	#sent	#words	vocab	len (hrs)
TACoS	cooking	26	127	360 sec	7,206	18,227	146,771	28,292	15.9
TACoS Multilevel	cooking	1	185	360 sec	14,105	52,593	2,000	-	27.1
MSVD	open	218	1970	10 sec	1,970	70,028	607,339	13,010	5.3
M-VAD	movie	-	92	6.2 sec	48,986	55,904	519,933	17,609	84.6
MPII-MD	movie	-	94	3.9 sec	68,337	68,375	653,467	24,549	73.6
MSR-VTT	open	20	7,180	20 sec	10,000	200,000	1,856,523	29,316	41.2
ActivityNet Captions	open	-	20,000	180 sec	-	100,000	1,348,000	-	849.0
ActivityNet Entities	social media	-	14,281	180 sec	52k	-	-	-	-
YouCook	cooking	6	88	-	Nil	2,688	42,457	2,711	2.3
YouCook II	cooking	89	2,000	316 sec	15.4k	15.4k	-	2,600	176.0
MP-II Cooking	cooking	65	44	600 sec	-	5,609	-	-	8.0
VideoStory	social media	-	20k	-	123k	123k	-	-	396.0
Charades	human	157	9,848	30 sec	-	27,847	-	-	82.01
VTW	open	-	18,100	90 sec	-	44,613	-	-	213.2
Kinetics 700	human	700	6,50,000	10 sec	700	-	-	-	1806

Table 2.2: Datasets and their characteristics

2.6 Evaluation Metrics

Text Generation is a tricky domain. Academics as well as the industry still struggle for relevant metrics for evaluation of the generative models' qualities. Every generative task is different, having its own subtleties and peculiarities. These metrics can be applied to the numerous tasks such as:

- Short or long-form text generation
- Machine Translation
- Summarisation
- Chatbots and dialogue systems
- Multimedia systems like speech2text, image/video captioning

Following are certain metrics that are used to evaluate the above natural language generation tasks:

2.6.1 BLEU: Bilingual Evaluation Understudy

BLEU [37] is a precision focused metric that calculates n-gram overlap of the target and generated texts. This n-gram overlap means that this evaluation scheme is word-position independent apart from n-grams' term associations. BLEU also consists of a brevity penalty i.e. a penalty applied when the generated text is too small compared to the target text.

2.6.2 ROUGE: Recall Oriented Understudy for Gisting Evaluation

ROUGE [38] is a set of metrics for evaluating automatic summarization of texts as well as machine translations. There are 3 types of ROUGE: ROUGE-N, the most common ROUGE type which means n-gram overlap. Second is ROUGE-L which checks for Longest Common Subsequence instead of n-gram overlap. The third is ROUGE-S which focuses on skip grams.

2.6.3 METEOR: Metric for Evaluation of Translation with Explicit Ordering

METEOR [39] is an metric that works on word alignments. It computes one to one mapping of words in generated and reference texts. Traditionally, it uses WordNet or porter stemmer. Finally, it computes an F-score based on these mappings. The metric was designed to fix some of the problems found in the more popular BLEU metric, and also produce good correlation with human judgement at the sentence or segment level. This differs from the BLEU metric in that BLEU seeks correlation at the corpus level.

2.6.4 CIDEr: Consensus based Image Description Evaluation

CIDEr [40] is a metric used to evaluate image descriptions that uses human consensus. It uses a triplet-based method of collecting human annotations to measure consensus followed by stemming and grouping the words into n-grams. This metric is especially useful in image and video annoated tasks.

2.6.5 WMD: Word Mover's Distance

WMD [41] is a fundamental technique for measuring the similarity of two documents. As the crux of WMD, it can take advantage of the underlying geometry of the word space by employing an optimal transport formulation. WMD leverages the results of advanced embedding techniques like word2vec and Glove. It suggests that distances between embedded word vectors are to some degree, semantically meaningful. It utilizes this property of word vector embeddings and treats text documents as a weighted point cloud of embedded words.

2.6.6 SPICE: Semantic Propositional Image Captioning Evaluation

SPICE [42] is an automated caption evaluation metric defined over scene graphs. It first transforms both generated and target captions into an intermediate representation that encodes semantic propositional content. It then creates a scene graph based on certain object classes, attribute types and relations which in turn is used to calculate the F-score b/w the generated and target captions.

2.7 Research Literature Gaps

2.7.1 Less Explored

The research on dense video captioning task is still on its early stage of development. There has been extensive work done in the field of computer vision on tasks like image classification, object localization and detection, image captioning, action classification on short video clips spanning few tens of seconds, etc. The computation costs in processing and cost of annotating long untrimmed videos is a major reason of less development in the field. Though research has been done on tasks like video summarization with a paragraph in the field of NLP, the models only utilizes the language features and completely discards the visual modality which can be very helpful in describing meaningful captions which cannot be captured by sole speech.

2.7.2 Use of Transformers

The existing research mostly utilizes CNN based models for feature extraction and proposal generation, given its performance legacy in image related tasks. The CNN based architectures require additional hand crafted components like anchor selection based on appearance and distribution of object to be predicted in the dataset,

threshold for non-max suppression to remove redundant predictions or multiple processing stages (eg. label generation) to detect arbitrarily shaped objects. Transformers on the other hand are more general architectures and have low inductive bias. Their structure has the capability to run on different modalities with very little modifications across modalities. With large amount of image and video data available across internet, transformer have the potential to become the state-of-the-art methods in computer vision due to their better generalizability.

2.7.3 Multimodal Architecture

The invention of Neural Networks is greatly inspired by the way human brain processes visual data and its ability to learn to recognize objects. Information in the real world usually comes as different modalities. Humans also utilize multiple modalities for processing and learning things. Majority of research work in DVC utilizes single modality of either just video frames or just speech. Thus it is important to design models that are able to jointly represent information from different modalities and make decision according to it.

2.7.4 Efficient Attention Mechanisms

Although there are many benefits of utilizing transformers, there are few concerns too. Transformers takes longer execution times for training and inference as compared to their CNN counterparts. The attention in transformer has quadratic complexity with respect to the number of input tokens. And the tokens in computer vision tasks are much larger than that in case of language tasks. These requires many decisions like input resolution and sliding stride size along both spatial and temporal dimension for feature generation. This will in turn affect the size of predicted objects or events. Hence methods are needed to reduce the execution time by smart selection of tokens for attention.

2.7.5 End-to-End Training

The video event segments and language are closely related and the caption information should be able to help localize temporal events in the video. Using independent or alternate training of modules disregards this relationship and predict inferior timestamps and captions. On the other hand, end-to-end training mechanisms allow proper gradient flow across modules and help in generating rich predictions.

Chapter 3

Methodology

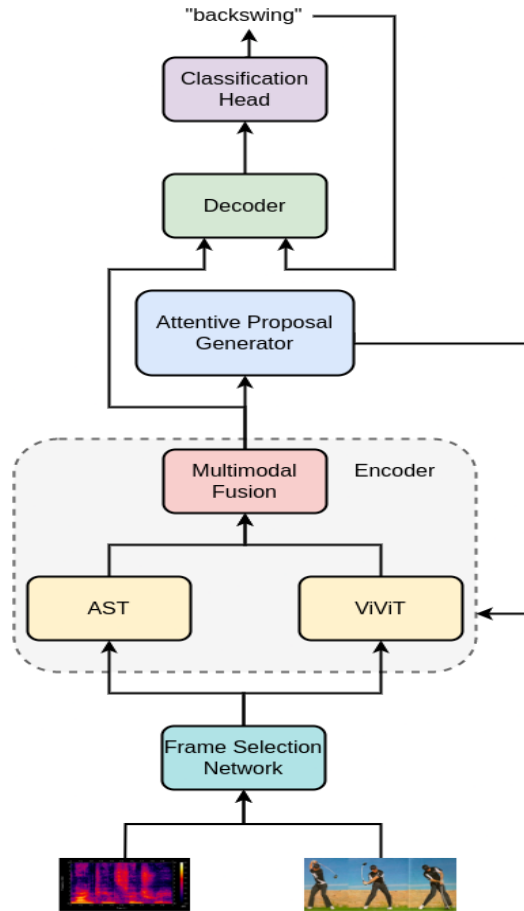


Figure 3.1: Proposed model for Dense Video Captioning

3.0.1 Multimodal Feature Extraction

Feature extraction is the backbone of our solution to tackle the task of Dense Video Captioning [2]. A rich feature space would enhance the representational power of the model, thereby leading to more meaningful and accurate event proposals. Most previous works have utilized one modality only (i.e. video) to generate feature vectors as input to the proposal generator. However, audio cues, in conjunction with video is a strong event indicator. These events are accompanied with a sharp change in their corresponding audio spectrogram which can be learnt by the model to better determine the precise boundary of events. Thus, we aim to combine

features generated using both video and audio. Currently, we are exploring two methods of fusion which include the cross-attention mechanism [6] and common space projection using contrastive learning [43].

Video Features

Video features are the most important part of the encoder. Without a robust and rich feature space for videos, the proposal generator would never be able to learn accurate event boundaries. Previous state-of-the-art video encoders [44], [45], [46] use CNN-based architectures. Although they have strong inductive bias and translational invariance, they fail to model long-range temporal dependencies which are of paramount importance when encoding videos. Moreover, CNNs require different architectures to model different modalities which can become complex when combining multiple modalities such as video and audio. Transformers [1] can overcome these barriers using its attention mechanism without compromising on its statistical and computational efficiency. Moreover, transformers can use the same building blocks across different modalities without many modifications. We aim to use the recently proposed ViViT model [47], a purely attention-based video encoder which has outperformed previous approaches across several datasets such as Kinetics 400 and 600, Epic Kitchens 100, Something-Something v2 and Moments in Time. Even though ViViT requires several orders of magnitude more training data as compared to its CNN counterparts, the authors of ViViT propose several methods to limit its training by using pretrained weights in conjunction with strong regularisation and specific fine-tuning.

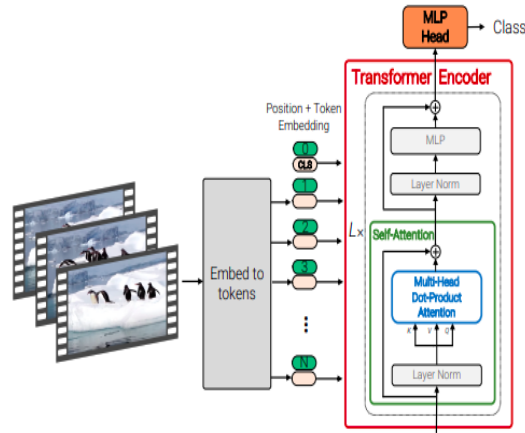


Figure 3.2: ViViT encoder by Dosovitskiy *et al.* (Image courtesy [47])

Audio Features

Audio is an important aspect of any video. Not only does audio suggest the duration of an event, it can also signify the magnitude or significance of that event. Thus, audio becomes a powerful accessory to images when representing a video. We aim to use the recently proposed Audio Spectrogram Transformer (AST) [48] which has outperformed current state-of-the-art models in audio classification. Moreover, AST uses pretrained weights from ViT [49], just like ViViT. We believe that this would thus, work cohesively with ViViT and lead to richer multimodal features.

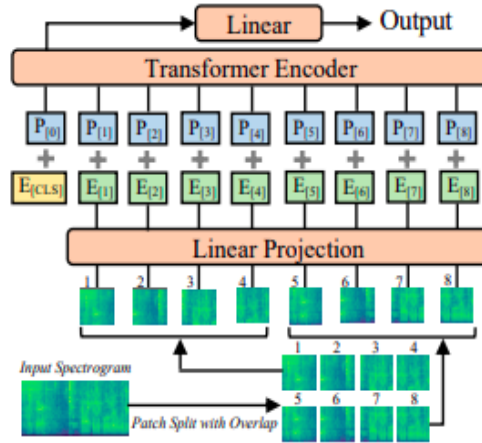


Figure 3.3: AST encoder by Gong *et al.* (Image courtesy [48])

3.0.2 Training

Knowledge Distillation

We aim to use the Knowledge Distillation framework (student-teacher paradigm) to train the model, either in individual modules or as a whole. For the video encoder (ViViT), we would use a strong CNN-based video classifier as the teacher model to introduce inductive bias within the transformer and reduce training. Touvron *et al* introduced a method for knowledge distillation [50] using a distillation token to compute the loss based on the softmax generated by the teacher model.

We also aim to train the proposal generator using a student-teacher paradigm with the current SOTA DVC models.

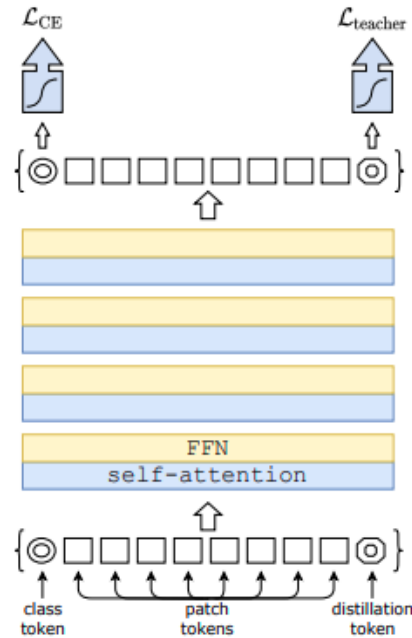


Figure 3.4: Distillation through attention by Touvron *et al.* (Image courtesy [50])

Chapter 4

Feature Extraction

For the task of Dense Video Captioning, a rich feature space would enhance the representational power of the model, thereby leading to more meaningful and accurate event proposals and captions. In most previous works, convolutional neural networks are commonly used to extract features from videos for any downstream tasks. Moreover, most localization methods use video features extracted by models that are trained for the task of Trimmed Action Classification (TAC), on datasets such as Kinetics[51] and Sports-1M, such as R(2+1)D[46], I3D[44], C3D, CSN [45] and others. Additionally, most previous works focus on using a single modality (i.e. video) or multiple modalities based only on video, like RGB and optical flow features, for example from I3D[44], to generate feature vectors as input to the downstream model. Some works have used multiple modalities including video (RGB, optical flow) and audio, like [6]. The motivation for using multiple modalities is so that the feature space can grow richer, and that multiple modalities can complement each other to give strong event indicators.

In this section, we explain how features are extracted from videos in the form of two modalities: video and audio. We use the *Video Vision Transformer (ViViT)* [47] for video feature extraction and the *Audio Spectrogram Transformer (AST)* [48] for audio feature extraction. We also consider the temporal sensitivity of features as a characteristic that can benefit DVC.

4.1 Feature Extraction Methodology

4.1.1 Video Features Extraction Using ViViT

Video features are the most important part of the encoder. Without a robust and rich feature space for videos, the downstream model would never be able to learn accurate event boundaries or captions. Previous state-of-the-art video encoders [44], [45], [46] use CNN-based architectures. Although they have strong inductive bias and translation invariance, they fail to model long-range temporal dependencies which are of importance when encoding videos. Moreover, CNNs require different architectures to model different modalities which can become complex when combining multiple modalities such as video and audio. Transformers [1] can overcome these barriers using their attention mechanism without compromising on statistical and computational efficiency. Moreover, transformers can use the same building blocks across different modalities without many modifications. We aim to use the recently proposed ViViT model [47], a purely attention-based video encoder which has outperformed previous approaches for action recognition such as Kinetics 400 and 600, Epic Kitchens 100, Something-Something v2 and Moments in Time. Even though ViViT requires several orders of magnitude more training data as compared to its CNN counterparts, the authors of ViViT propose several methods to limit its training requirements by using pretrained weights in conjunction with strong regularization and specific fine-tuning.

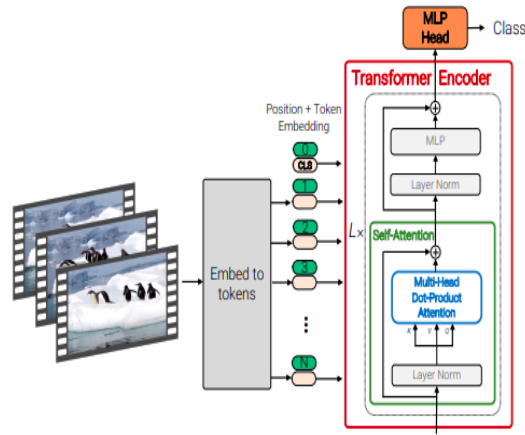


Figure 4.1: ViViT encoder by Dosovitskiy *et al.* (Image courtesy [47])

Arnab *et al* propose two methods for mapping a video to its corresponding embeddings for input to ViViT, *Uniform frame sampling* and *Tubelet embedding* in [47]. We use the uniform frame sampling method in our feature extraction framework. Arnab *et al* also propose different attention mechanisms for the ViViT model; we use the *spatio-temporal attention* mechanism. The flow of data is as follows:

1. Consider the original video clip X as consisting of T frames, each of width W and height H . The video has C channels. Hence the shape of video frames of X is (T, H, W, C) .
2. We perform certain preprocessing steps on the video frames:
 - (a) We resize the spatial dimensions of the video frames. The shorter edge is resized to 256, and the longer edge is resized maintaining the aspect ratio.
 - (b) We normalize all channels of the video frames.
 - (c) Generalizability and robustness of the features motivate data augmentation of the video frames. To this end, we randomly flip the video frames horizontally, and take random crops of spatial dimensions 224×224 during training. Note that for validation and inference, we skip the random horizontal flip step and take a center crop of the same dimensions.
3. Now, the uniform frame sampling method proposed by Arnab *et al* in [47] is used to extract tokens mapped to a 768-dimensional space. This involves using a three-dimensional convolutional layer to project the frames to patch embeddings, so that the shape of the embeddings are $(T, N, 768)$, where N is the number of patches.
4. The patch embeddings are then flattened to $(T \times N, 768)$, and concatenated with a *class token* for input to the encoder layers of ViViT.
5. For each encoder layer, the input is subject to (1) positional encoding, (2) layer normalization, (3) multi-head self-attention across all tokens, (4) layer normalization and (5) two fully connected layers with GeLU activation and dropout.
6. Finally, the class token from the output of the last encoder layer is used as the feature representation of the clip X .

More details about the different frame sampling methods, attention mechanisms and the ViViT model are given in B.4

4.1.2 Audio Features Extraction using AST

Audio is an important aspect of any video. Not only does audio suggest the duration of an event, it can also signify the magnitude or significance of that event. Thus, audio becomes a powerful accessory to images when representing a video. Convolutional neural networks (CNNs) have been widely used to learn representations from raw spectrograms for end-to-end modelling, as the inductive biases inherent to CNNs such as spatial locality and translation equivariance are believed to be helpful. In order to better capture long-range global context, a trend is to add a self-attention mechanism on top of the CNN. Such CNN-attention hybrid models have achieved state-of-the-art (SOTA) results. We aim to use the recently proposed Audio Spectrogram Transformer (AST) [48] which has outperformed current state-of-the-art models in audio classification. Audio Spectrogram Transformer is a convolution-free, purely attention based model that is directly applied to an audio spectrogram and can capture long-range global context even in the lowest layers.

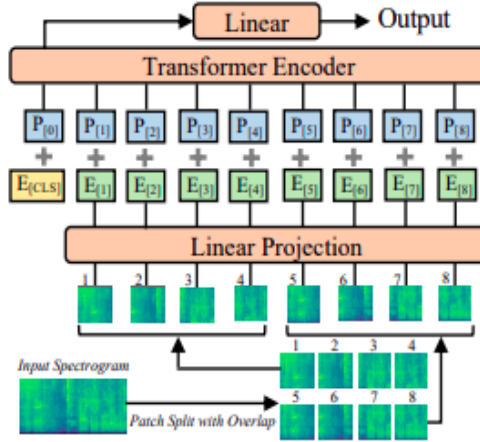


Figure 4.2: AST encoder by Gong *et al.* (Image courtesy [48])

The flow of data is as follows:

1. The audio signal is first pre-emphasised, then sliced into (overlapping) frames and a window function is applied to each frame. Following that, we perform a Short-Time Fourier Transform on each frame, calculate the power spectrum, and then compute the filter banks.
2. Thus a sequence of 128-dimensional log Mel filterbank (fbank) features is computed with a 25ms Hamming window every 10ms from the input audio waveform of t seconds.
3. The spectrogram is then split into a sequence of N 16×16 patches with a time and frequency overlap of 6, where N is the number of patches and the Transformer's effective input sequence length.
4. Each 16×16 patch is then flattened into a 1D patch embedding of size 768 using a linear projection layer.
5. Trainable positional embedding is added to each patch embedding which allows the model to capture the spatial structure of the 2D audio spectrogram.
6. The class token is then concatenated at the beginning of the sequence, which is subsequently fed into the encoder layer of AST.
7. Finally, the class token from the output of the last encoder layer is used as the feature representation of the clip.

4.2 Temporal Sensitivity of Features

As mentioned previously, most localization methods use video features extracted by models that are trained for Trimmed Action Classification (TAC) tasks. These features are not necessarily optimal for Temporal Action

Localization (TAL) and other dependent tasks such as Action Proposal Generation and Dense Video Captioning. Alwassel *et al*, in their paper titled *TSP: Temporally Sensitive Pretraining of Video Encoders for Localization Tasks*, claim that this is due to the inherent lack of temporal sensitivity in features from TAC trained video encoders [35]. Yet, these features are prominently used for TAL and DVC tasks. This discrepancy is explained by the following reasons:

- Due to a lot of existing research for TAC tasks, a lot of established models exist that can be used off the shelf for video feature extraction
- The video encoders usually cannot be directly trained along with downstream tasks of TAL or DVC due to resource constraints; it is impractical to fit long untrimmed videos in commodity GPUs without drastically downsampling them in terms of space or time.

To solve this problem, Alwassel *et al* introduced a supervised pre-training paradigm for TAL that aims to instill the ability to produce temporally sensitive features; video encoders are trained to explicitly discriminate between foreground and background clips in untrimmed videos. They show that TSP improves performance for TAL, action proposal and DVC tasks consistently on different datasets as well as with different downstream predictive models and architectures, clearly establishing that temporally sensitive features contribute positively to the task at hand. TSP involves training encoders with a downstream task involving two classification heads: (1) *action classification* head and (2) *temporal region classification* head. This architecture is applied to clips from untrimmed videos. Optionally, along with local features of every clip, a *global video feature (GVF)* is also used for the temporal region classification to further improve temporal sensitivity [35]. A summary of their work is given in B.6.

4.2.1 Temporally Sensitive Pretraining of Proposed Feature Extractors

We extend the TSP framework proposed by Alwassel *et al* [35] to use temporally sensitive features in our work for DVC.

Architecture

The model used for TSP consists of one or more *backbone* architectures for particular modalities, which in our case are ViViT - for video modality, and AST - for audio modality. The constructed model is such that different backbones and their modalities can be plugged in and out as required for experimentation. The input modality data of a clip is fed in to each backbone separately, and the output feature representations from each backbone are combined and passed through a fully connected layer to give a final feature representation. This feature representation is fed into the action classification head and the temporal region classification head to give the respective outputs. Note that we do not use a GVF in the current implementation of our pretraining framework, and plan to include this in later work.

Data Flow

In our case, the following data flow is followed (considering two backbones: ViViT and AST):

1. Video frames are extracted from the clips, and are preprocessed and augmented as mentioned in 4.1.1 making the features at this level of the shape (B, C, T, H, W) , where B is the batch size, C is the number of channels (3 for RGB), T is the temporal length of the clip in terms of number of frames and H and W are the height and width (spacial dimensions) of each frame in the clip respectively.
2. Audio frames are extracted from the clips as explained in 4.1.2. At this stage, they are of the shape (B, N, F) , where N is the number of time frames and F is the frequency bins.
3. Video features of the shape (B, C, T, H, W) are fed into ViViT, and the class token of shape $(d_v,)$ of ViViT is used as the feature representation f_v of the video modality of the clip.
4. Audio features of the shape (B, N, F) are fed into AST, and the class token of the shape $(d_a,)$ of AST is used as the feature representation f_a of the audio modality of the clip.

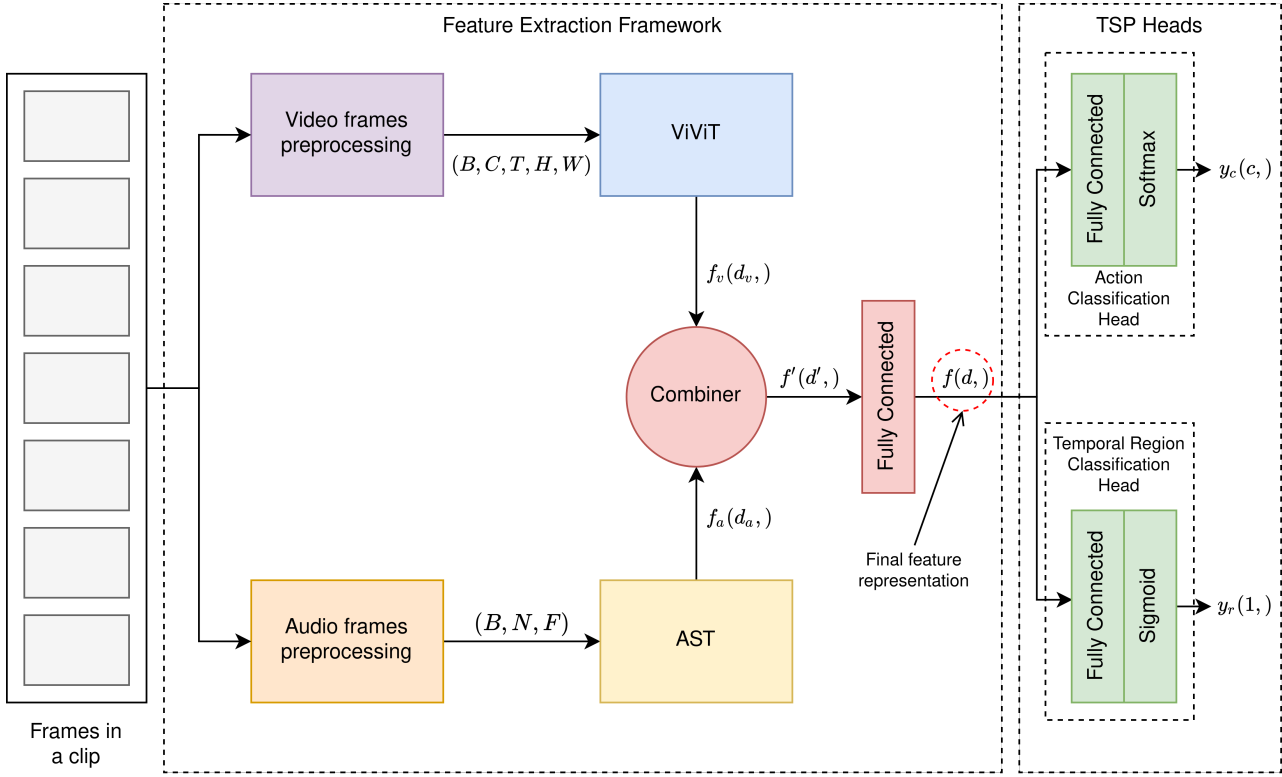


Figure 4.3: The proposed feature extraction framework with ViViT and AST backbones. TSP heads for action classification and temporal region classification are shown on the right.

5. These two features, f_a and f_v , are then combined using some function, for example addition or concatenation. In general, if there are n backbones and we have features $f_1, f_2, f_3, \dots, f_n$, this combiner function must combine all these and give a single-dimensional feature f' of dimensions $(d',)$.
6. The feature f' is fed into a fully connected layer that converts it into a feature f of shape $(d,)$. This feature f is the final feature representation that is used for downstream tasks.
7. f is fed to the action classification head, which consists of a simple fully connected layer followed by softmax activation layer, to give the action classification logits y_c as output.
8. f is also fed (optionally combined with the GVF, f_g , to improve the temporal sensitivity of the feature representation) to the temporal region classification head, which consists of a fully connected layer followed by the sigmoid activation layer, to give the binary temporal region classification logit y_r as output.

The global video feature is extracted in the same manner as local features, with the additional step of *pooling* across all local features f_1, f_2, \dots, f_n for a video.

Loss functions

As mentioned, there are two classification heads on the top of the TSP framework:

- **Action Classification:** This head predicts action classes for a given input clip. The head can be represented as $f \rightarrow y_c$. The cross-entropy loss function is used for this head, with mean reduction (L_c):

$$L_{CE} = -\frac{1}{n} \sum_{i=1}^n \mathbf{y} \log(\hat{\mathbf{y}})$$

- **Temporal Region Classification:** This head predicts temporal region classes for a given input clip, i.e. whether the particular clip lies in an action (foreground) temporal region or non-action (background) temporal region. The head can be represented as $f \rightarrow y_r$. For binary classification, the cross entropy function reduces to binary cross entropy (L_r).

The two losses, L_c and L_r are given loss coefficients α_c and α_r for relative importance for contributing towards the total loss:

$$L = \alpha_c L_c + \alpha_r L_r$$

We evaluate the effectiveness of this temporally sensitive pretraining in the Results section.

Appendix A

Supplementary Details

Appendix B

Paper Summaries

B.1 Dense-Captioning Events in Videos

B.1.1 Overview

Ranjay Krishna *et al*, in their 2017 paper titled *Dense-Captioning Events in Videos* [2], proposed the problem of *dense video captioning*. The network was able to output overlapping events and events with varying time scales. The captioning module output also considered past and future context of the representative event. They also introduced new dataset with over 20K annotated videos for this problem - ActivityNet Captions.

B.1.2 Datasets

- ActivityNet Captions

B.1.3 Performance

Krishna *et al* compared captioning results with LSTM-YT[52], S2VT[53] and H-RNN[54] using BLEU@1-4, METEOR and CIDEr. Since these models were only for captioning, the results were compared by providing ground-truth proposal boundaries. The metrics showed better results even when their pipeline generated captions on learnt proposals. The model was tested with different type of contexts input in captioning model, and full-context (*both past and future*) proved to be better for longer videos. The results of event localization depicted how multi-stride network performed better for videos with more number of events. The model also reported state-of-the-art results for the task of video and paragraph retrieval

B.1.4 Methodology

Krishna *et al* architecture consists of two main components:

1. Event proposal module
2. Captioning module with context

Flow of pipeline

1. C3D features of sampled video frames.
2. Sampling done at different strides to capture all kind of events.
3. Features fed into LSTM type of network inspired by DAPs[18] for generating proposals.
4. Captioning module used a LSTM network with following inputs:

- (a) hidden representation of past events h_i^{past}

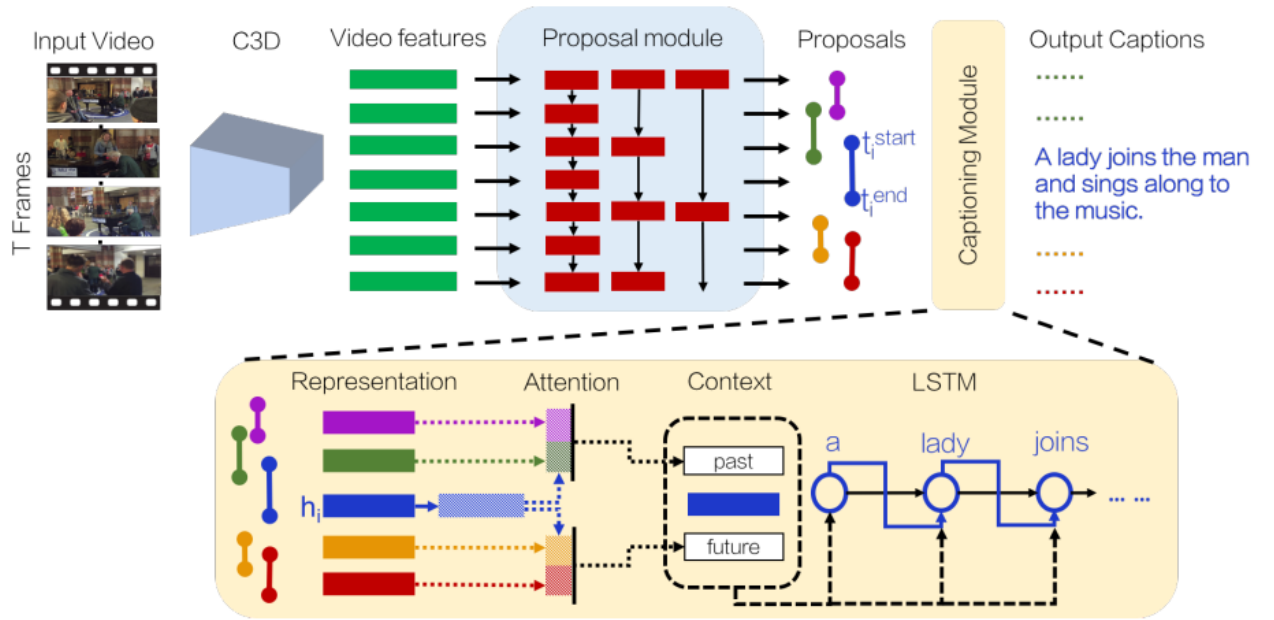


Figure B.1: Pipeline introduced by Krishna *et al* (Image courtesy [2])

- (b) hidden representation of current event h_i
- (c) hidden representation of future events h_i^{future}

B.1.5 Conclusion

Krishna *et al* tackled the task of captioning videos with correlation between events. They also showed the results on streaming videos with current event relating to only past events.

B.2 A Better Use of Audio-Visual Cues: Dense Video Captioning with Bi-modal Transformer

B.2.1 Overview

Iashin and Rahtu, in their 2020 paper titled *A Better Use of Audio-Visual Cues* [6], approached the task of Dense Video captioning by using two types of input features: video as well as audio. They identified the lack of research into utilizing the audio track of videos to generate dense captions, despite the natural co-occurrence of the two. They introduce a novel *Bimodal Transformer* which has the similar encoder-decoder architecture as the traditional Transformer model, but generalizes it to two input modalities: audio features (using VGGish [55]) and visual features (using I3D [44]). Their framework also consists of a *multiheaded bimodal proposal generator*, inspired by the YOLO object detector [30] which generates event proposals for caption generation.

B.2.2 Datasets

- ActivityNet Captions [2]

B.2.3 Performance

Iashin *et al* report state-of-the-art-performance for the dense video captioning task on the published validation subset of the ActivityNet Captions dataset in terms of BLEU@3-4 metrics and near state-of-the-art performance in terms of the METEOR metric. Considering the task of proposal generation in isolation, they report state-of-the-art performance with the F1 score of 60.27%, on the same dataset.

B.2.4 Methodology

Iashin *et al* use an Inflated 3D Network pre-trained on the *Kinetics* dataset [51] for extracting visual features from video frames. For audio features, the VGGish network [55] is used. These features are stacked in sequence for each video. The caption tokens are embedded with pre-trained GloVe [56].

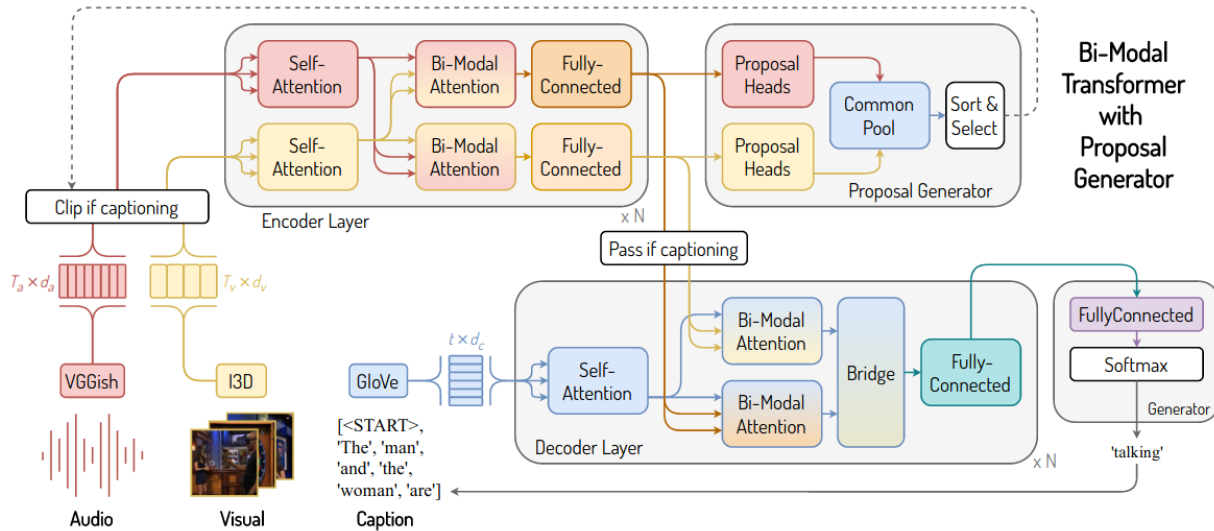


Figure B.2: Dense Video Captioning Framework introduced by Iashin *et al* (Image courtesy [6])

The design proposed by Iashin *et al* consists of two major functional components:

- Bi-modal Transformer
- Multiheaded Proposal Generator

For inference, the features flow through the components as follows:

1. Both feature sequences, audio A and visual V are self-attended and then passed through an N -layered encoder to produce bi-modal sequence representations using novel bi-modal multi-headed attention blocks to fuse features from both sequences. These bimodal attention blocks output two sequences of features, A^v (visual attended audio features) and V^a (audio attended visual features).
2. A^v and V^a are then given to the Proposal Generator, which consists of two distinct sets of *proposal generator heads*, one set for each input modality. Each proposal generator head makes proposal predictions at each time step independently, forming a common pool of cross-modal predictions. Top k predictions are selected (using a confidence score) from this pool to be the output of the proposal generator: the proposals.
3. The generated proposals are used to clip the input feature sequences which are then given to the Bimodal Encoder. The same self-attention followed by bimodal attention is carried out on these clipped features.
4. The encoder's outputs are passed to the bi-modal attention blocks in every decoder layer of the Bimodal Decoder, along with the representation of previously generated caption words. The caption embeddings are self-attended, and then carries out bi-modal encoder-decoder attention, producing $C_n^{A^v}$ (audio-visual attended previous captions) and $C_n^{V^a}$ (visual-audio attended previous captions). After a bridge connection and positional encoding (to add order information to the permutation invariant transformer architecture), two fully connected layers with ReLU activation between them and a softmax activation in the final layer are used to model the distribution of the next caption token over the vocabulary.

Proposal Generation Module

The design of individual proposal generator heads in the Proposal Generator is inspired by the YOLO object detector [30]. A set of anchors is determined apriori by running K-Means Clustering on the ground truth event lengths. Each centroid of the cluster is taken as an anchor in the anchor set ψ . The distance metric used here is the Euclidean distance, in contrast to IoU (Intersection over Union) used in YOLO. Each proposal generator head predicts three values: (1) center of the event c , (2) log coefficient l and (3) objectness score o . Using these values, the temporal boundaries and confidence of the event predicted are calculated:

$$\begin{aligned} center &= p + \sigma(c) \\ length &= anchor_length \cdot e^l \\ confidence &= \sigma(o) \end{aligned}$$

Each proposal generator head consists of three 1D convolutional layers, through which the sequence length is preserved using padding and unit stride. The first convolutional layer has kernel size from a predetermined set of kernel sizes (calculated for each modality using K-Means clustering on event ground truth event lengths; the idea is to match receptive field sizes with event sizes). The other two layers have unit kernel size.

The loss for center c and log coefficient l , collectively called *Localization error* is the Mean Squared Error (MSE), while the confidence loss is a weighted sum of Binary Cross Entropy (BCE) of objectness and no-objectness loss.

Captioning Module

The captioning module is essentially the Bimodal Transformer, consisting of a bimodal encoder and a bimodal decoder. The input to the bimodal encoder during captioning is feature sequences A and V , which temporally correspond to a proposal. These two sequences go through N encoder layers, each of which consist of (1) multiheaded self-attention of the two input sequences, (2) multiheaded bi-modal attention of the two sequences, giving A^v and V^a and then (3) two position-wise fully-connected layers with residual connections. These sublayers have separate trainable weights for both modalities.

The bimodal decoder is given the previous sequence of caption tokens as input, along with the output of the encoder A_N^v and V_N^a . Each decoder layer consists of the following sublayers: (1) self-attention of caption

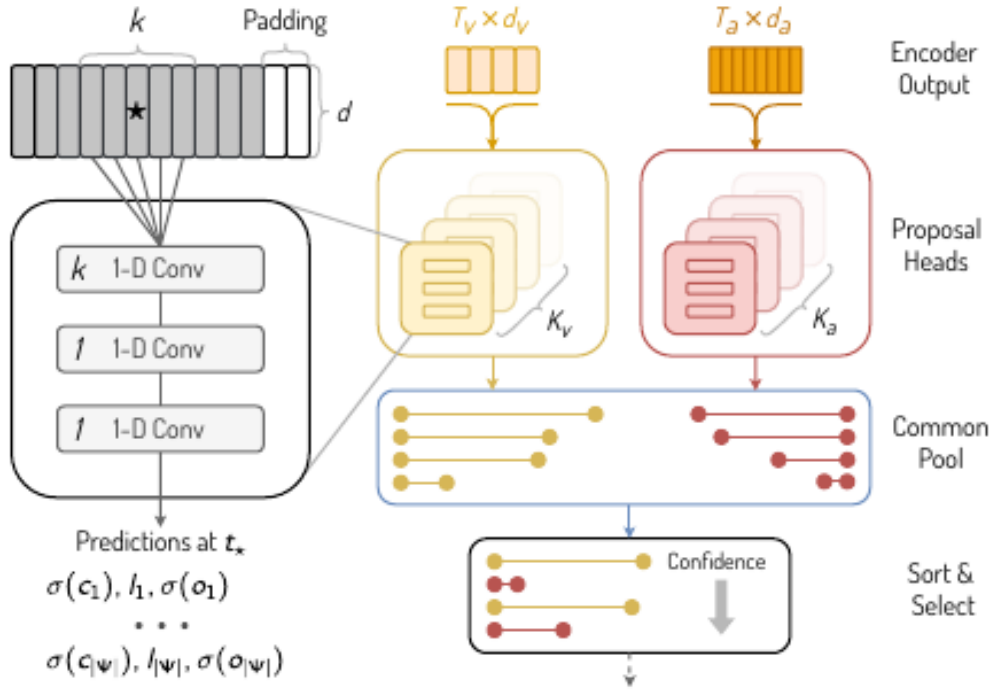


Figure B.3: Bi-modal Multiheaded Proposal Generator: Architecture (Image courtesy [6])

tokens, (2) bimodal attention of captions with A^v and V^a , (3) a bridge connection to concatenate and combine $C_n^{A^v}$ (audio-visual attended previous captions) and $C_n^{V^a}$ (visual-audio attended previous captions) and (4) two position-wise fully connected layers which act as the classification layers for the next caption token.

Training Procedure

First, the captioning module is trained using ground truth proposals and their given captions using KL-divergence loss and label smoothing. The bimodal encoder weights are then frozen and then the proposal generator is trained using the trained bimodal encoder. Each proposal head uses Mean Squared Error for localization and binary cross-entropy for proposal loss.

B.2.5 Conclusion

The work by Iashin *et al* provides a compelling case for exploring the use of multiple modalities for the task of Dense Video Captioning. Their ablation study shows that performance with using only visual features is more than using only audio features, indicating visual modality has a stronger signal for video understanding than audio modality; however, performance using both modalities gives consistently better results than single modalities in all settings.

B.3 Vision Transformer (ViT) - An image is worth 16x16 words: Transformers for image recognition at scale

B.3.1 Overview

Dosovitskiy *et al*, in their 2020 paper titled *An image is worth 16x16 words: Transformers for image recognition at scale* [49], aim to replace CNN based architectures with transformers for a multitude of computer vision tasks. They propose a transformer-based architecture inspired by Vaswani *et al* [1], replacing the word embeddings with non-overlapping ‘patch embeddings’, extracted from images. The authors show that if ViT is pre-trained on large amounts of data, it can produce better results as compared to its CNN counterparts in various image recognition benchmarks.

B.3.2 Datasets

- Pre-trained on:
 - ILSVRC-2012 ImageNet dataset with 1k classes and 1.3M images
 - ImageNet 21k dataset with 21k classes and 14M images
 - JFT with 18k classes and 303M high-resolution images
- Benchmark tasks
 - ImageNet (validation labels)
 - CIFAR 10/100
 - Oxford-IIIT Pets
 - Oxford Flowers-102
 - 19-task Visual Task Adaptation Benchmark (VTAB) classification suite.

B.3.3 Performance

Dosovitskiy *et al* report state-of-the-art-performance for various image recognition tasks (datasets mentioned above), when compared with current state-of-the-art models such as Big Transfer (BiT) [57] and Noisy Student [58].

B.3.4 Methodology

Dosovitskiy *et al* use the original transformer model [1], but with different input embeddings that are extracted from images. The transformer model then outputs a classification ([class]) token which represents the entire image. This [class] token is then used for downstream tasks (mentioned above).

Input: Patch embeddings

The input consists of flattened $n \times n$ patches of an image. These patches are non-overlapping and span the entire image.

- The original image $x = R^{H \times W \times C}$.
- This is converted to a flattened 2D patch $x_p = R^{N \times (P^2 \cdot C)}$ where $H \times W$ is the resolution of the image, C is the number of channels, $P \times P$ is the resolution of each patch.
- $N = HW/P^2$ represents the number of patches extracted from the image and serve as the sequence length for the transformer. These are the patch embeddings.
- Position embeddings are added to the above patch embeddings. These can be:
 - 1D, like in transformers.

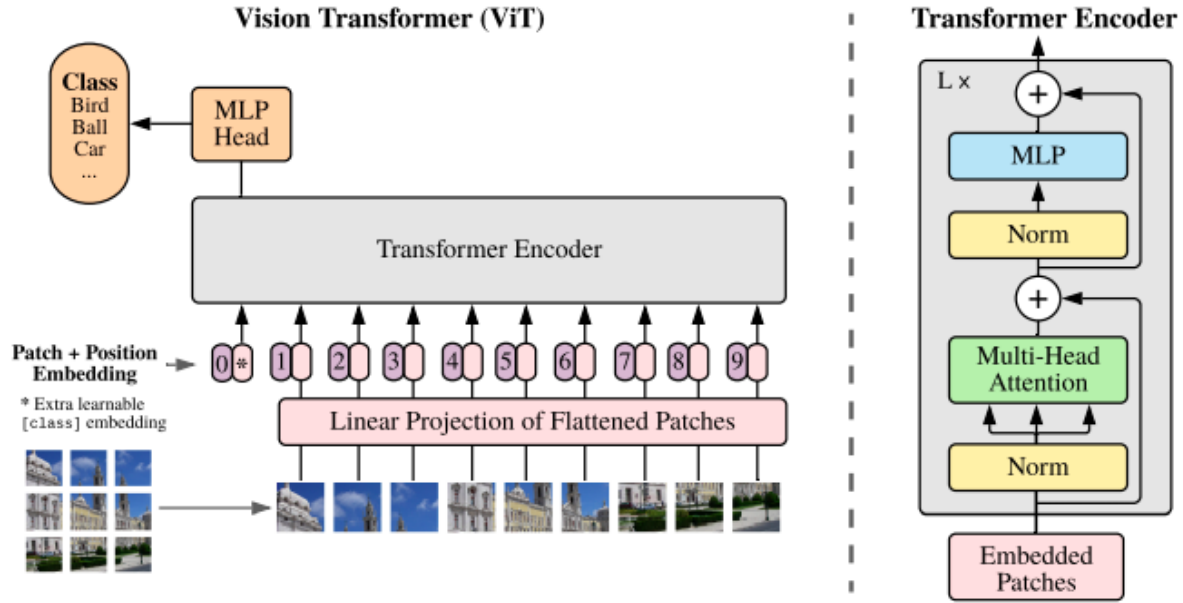


Figure B.4: Video Transformer introduced by Dosovitskiy *et al* (Image courtesy [49])

$$\begin{aligned}
 \mathbf{z}_0 &= [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}}, & \mathbf{E} &\in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D} \\
 \mathbf{z}'_\ell &= \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, & \ell &= 1 \dots L \\
 \mathbf{z}_\ell &= \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, & \ell &= 1 \dots L \\
 \mathbf{y} &= \text{LN}(\mathbf{z}_L^0)
 \end{aligned}$$

Figure B.5: Video Transformer introduced by Dosovitskiy *et al* (Courtesy [49])

- 2D, based on row and column indices.
- Relative distances between patches to capture the true spatial positioning of each patch within the image.
- A learnable [class] token is prepended to the image sequence like the BERT model. \mathbf{z}_0 is the input [class] token and \mathbf{z}_L is the output token. \mathbf{z}_L is a vector that represents the entire image and can be used for downstream tasks.
- These embeddings are mapped onto a D-dimensional space which is used by the encoder (ViT).

Transformer Encoder

The transformer model proposed in [1] is inherently the same one used by ViT. The only difference is that a normalisation layer is added before the multi-headed self attention block. Additionally, ViT used the GELU non-linearity as part of the multi-layer perceptron block.

Training Procedure

Pretrained on datasets such as ImageNet, ImageNet-21k, JFT-300M for better performance.

- Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, batch size = 4096, weight decay = 0.1

- Strong regularization is used throughout the encoder. Dropout, when used, is applied after every dense layer except for the the query-key-value projection layers and directly after the addition of positional embeddings
- Training is done on a resolution of 224(14×14)

B.3.5 Conclusion

The work by Dosovitskiy *et al* provides a compelling case for using transformers instead of CNN-architectures in image recognition tasks. Even though it requires large amount of pre-training, it can be scaled much better than its CNN counterparts and can meet the performance (and even exceed it in certain cases) of current state-of-the-art CNN-based models.

B.4 ViViT: A Video Vision Transformer

B.4.1 Overview

Arnab *et al*, in their 2021 paper titled *ViViT: A Video Vision Transformer* [47], extend the Vision Transformer (ViT) [49] to tackle video recognition tasks. They propose various attention architectures to model the spatial and temporal interaction within video frames. They also propose embedding extraction methods for videos which are on similar lines to those done for images.

B.4.2 Datasets

- Pre-trained on:
 - ILSVRC-2012 ImageNet dataset with 1k classes and 1.3M images
 - ImageNet 21k dataset with 21k classes and 14M images
 - JFT with 18k classes and 303M high-resolution images
- Benchmark tasks
 - Kinetics 400/600
 - Epic Kitchens
 - Something-Something v2
 - Moments in Time

B.4.3 Performance

Arnab *et al* report state-of-the-art performance for various video recognition tasks (datasets mentioned above), when compared with current state-of-the-art models.

B.4.4 Methodology

Arnab *et al* use the ViT model [49] as its foundation, but with different input embeddings that are extracted from video and different attention blocks. The transformer model then outputs a classification ([class]) token which represents the entire video. This [class] token is then used for downstream tasks (mentioned above).

Input: Token embeddings

Two methods for mapping a video to its corresponding embeddings are proposed, namely:

- *Uniform frame sampling*
 - The original video $x = R^{T \times H \times W \times C}$.
 - From this, we extract $N = n_t \cdot n_h \cdot n_w$ tokens which are subsequently mapped to d-dimensional space. Here, n_t corresponds to the number of frames in the video.
- *Tubelet embedding*
 - This method models the 3D convolutions or "tubes".
 - $n_t = T/t$ frames are sampled, each with $n_h \cdot n_w$ non-overlapping patches.
 - From this, we extract $N = n_t \cdot n_h \cdot n_w$ tokens which are subsequently mapped to d-dimensional space.
- The final output consists of $x = R^{N \times d}$
- A learnable [class] token may be prepended to the token sequence based on the attention architecture being used.

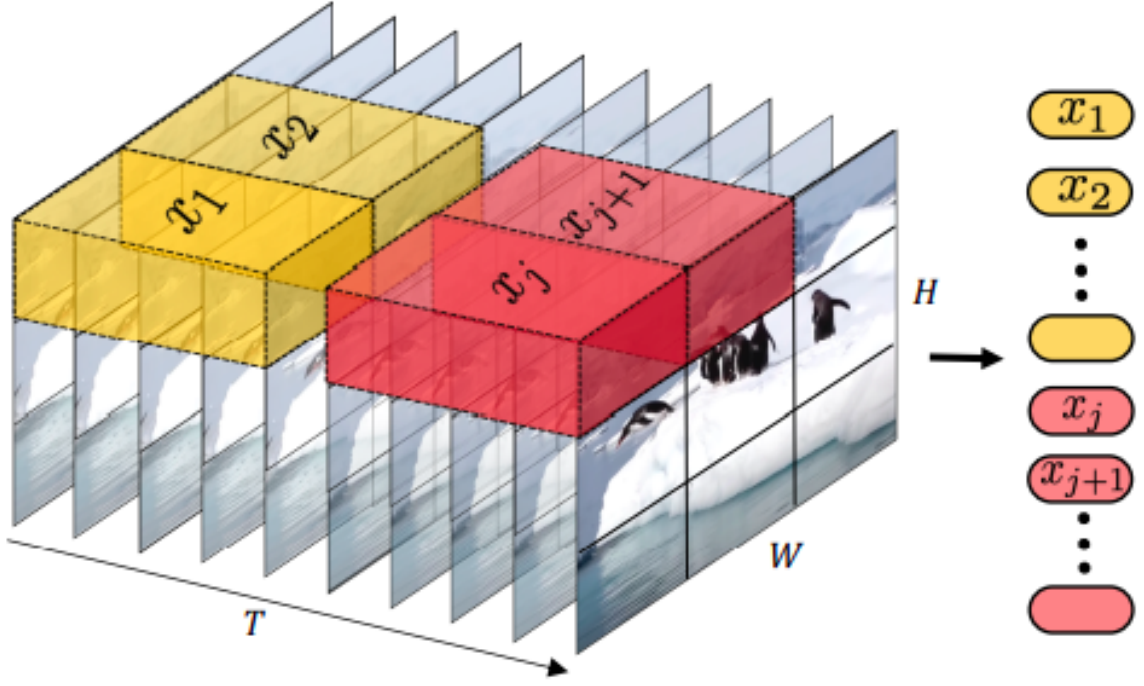


Figure B.6: Video Transformer introduced by Arnab *et al* (Courtesy [47])

Transformer Encoder

ViViT extends the transformer model proposed in [49] by making use of different attention architectures to model the spatio-temporal interaction between the token embeddings.

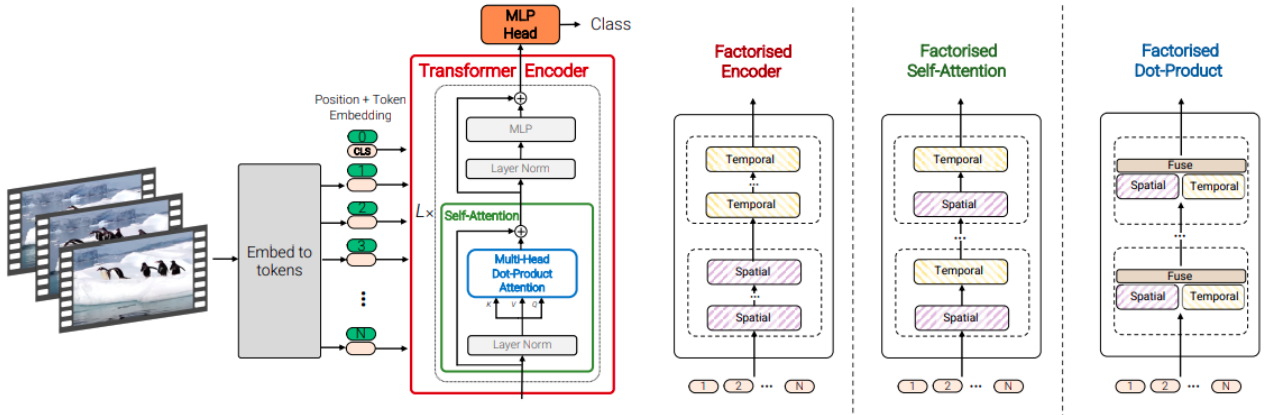


Figure B.7: Video Transformer introduced by Arnab *et al* (Image courtesy [47])

Four attention architectures are proposed, namely:

- *Spatio-temporal attention*
 - This architecture mimics the multi-headed self attention block [1] across all $N \times d$ tokens.
 - However, it has quadratic time complexity as it models all pairwise interactions between all spatio-temporal tokens.
- *Factorised encoder*

- This architecture consists of two transformer encoders to model spatial and temporal interactions between the tokens separately.
- First, the spatial encoder gets $n_h \cdot n_w$ tokens in order to attend to tokens within the same temporal index (same frame/tube).
- Next, n_t [class] tokens representing $n_h \cdot n_w$ tokens of each frame are given as input the the temporal encoder.
- The temporal encoder outputs another [class] token which is the collective representation of all the frames in the video.
- Model 2 does have more parameters than model 1, but consists of fewer floating point operations due to a decrease in the number of tokens per encoder.
- *Factorised Self-Attention*
 - This architecture consist of two multi-headed self attention blocks within a single encoder.
 - The first one models spatial interactions between tokens and the second one, temporal.
 - The first attention block gets tokens $x = R^{n_t \times n_h \cdot n_w \times d}$ which include tokens within the same temporal index.
 - The second attention block gets tokens $y = R^{n_h \cdot n_w \times n_t \times d}$ which include tokens within the same spatial index.
 - Model 3 has fewer parameters but the same complexity as model 2.
- *Factorised Dot-Product Attention*
 - This architecture aims to fuse the spatial and temporal attention within the same multi-headed self attention block, but with different heads.
 - The spatial attention heads get $x = R^{n_t \times n_h \cdot n_w \times d}$ tokens as input whereas the temporal attention heads get $y = R^{n_h \cdot n_w \times n_t \times d}$ as input.
 - Model 4 has the same number of parameters as model 1 and the same complexity as models 2 and 3.
- The final output can either be a [class] token or a 1D vector obtained by average pooling all output vectors.

Training Procedure

Pre-trained weights of the ViT model are leveraged due to the size of various video recognition tasks (many orders of magnitude smaller than image recognition datasets). The model can then be trained on various video recognition datasets such as Kinetics and Epic Kitchens.

- Kinetics
- Strong regularization is used throughout the encoder. Dropout, when used, is applied after every dense layer except for the the query-key-value projection layers and directly after the addition of positional embeddings
- Training is done on a resolution of 224(14×14)

B.4.5 Conclusion

The work by Arnab *et al* provides a compelling case for using transformers instead of CNN-architectures in video recognition tasks. Even though it requires large amount of pre-training, it can be scaled much better than its CNN counterparts and can meet the performance (and even exceed it in certain cases) of current state-of-the-art CNN-based models.

B.5 AST: Audio Spectrogram Transformer

B.5.1 Overview

Gong *et al* in their 2021 paper titled *AST: Audio Spectrogram Transformer* [48] proposed a convolution-free, solely attention-based model that can capture long-range global context even in the lowest layers and is directly applied to an audio spectrogram. They also suggest a method for transferring knowledge from the ImageNet-pretrained Vision Transformer (ViT) to AST, which can greatly enhance performance.

B.5.2 Datasets

- AudioSet
- ESC-50
- Speech Commands

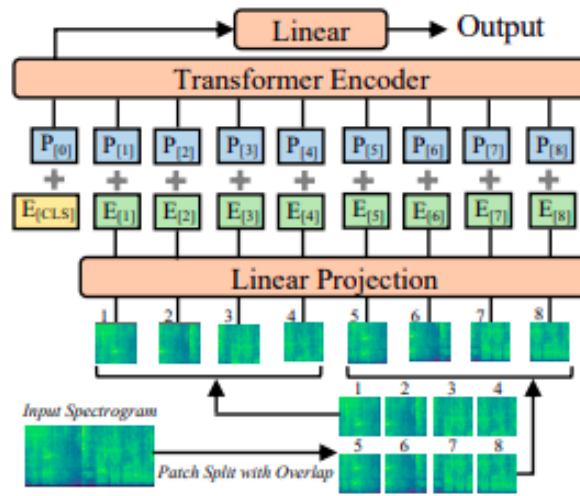


Figure B.8: Audio Spectrogram Transformer by Gong *et al* (Courtesy [48])

B.5.3 Methodology

- A sequence of 128-dimensional log Mel filterbank (fbank) features is computed with a 25ms Hamming window every 10ms from the input audio waveform of t seconds. As a result, the AST receives a $128 \times 100t$ spectrogram as input.
- The spectrogram is then split into a sequence of N 16×16 patches with a time and frequency overlap of 6, where N is the number of patches and the Transformer's effective input sequence length.
- Each 16×16 patch is then flattened into a 1D patch embedding of size 768 using a linear projection layer. This linear projection layer is called patch embedding layer.
- Since the Transformer architecture does not capture the input order information and the patch sequence is also not in temporal order, trainable positional embedding (also of size 768) is added to each patch embedding which allows the model to capture the spatial structure of the 2D audio spectrogram.
- The CLS token is then appended to the beginning of the sequence, which is subsequently fed into the Transformer.
- A Transformer consists of several encoder and decoder layers. For classification tasks, AST only uses the encoder of the Transformer.

- The [CLS] token output from the Transformer encoder serves as the audio spectrogram representation. The audio spectrogram representation is mapped to labels for classification using a linear layer with sigmoid activation.
- AST model also allows cross-modality transfer learning since images and audio spectrograms have similar formats. So ImageNet-pretrained CNN weights are used as initial CNN weights for audio classification training.

B.5.4 Training procedure

- Model is pre-trained with the ImageNet dataset.
- Model is trained with a batch size of 12, the Adam optimizer, and uses binary cross-entropy loss.
- For balanced set experiments, the initial learning rate of $5e-5$ is used and the model is trained for 25 epochs, the learning rate is cut into half every 5 epochs after the 10th epoch.
- For full set experiments, the initial learning rate of $1e-5$ is used and the model is trained for 5 epochs, the learning rate is cut into half every epoch after the 2nd epoch.
- Mean Average Precision (mAP) is used as the main evaluation metric.

B.5.5 Conclusion

Gong et al present Audio Spectrogram Transformer which is a convolution-free, purely attention-based model for audio classification which features a simple architecture and superior performance.

B.6 Temporally Sensitive Pretraining of Video Encoders for Localization Tasks

B.6.1 Overview

Most localization methods use video features extracted by models that are trained for Trimmed Action Classification (TAC). They’re not necessarily suited for Temporal Action Localization (TAL), since TAC-pretrained features tend to be temporally insensitive. For example, R(2+1)D, I3D and C3D have become the de facto video feature extractors for TAL, Action Proposal and DVC tasks; these are trained on TAC. The reasons for most work using TAC trained video encoders are (1) many established models exist for TAC and (2) it is impractical to fit untrimmed videos in commodity GPUs without drastically downsampling space or time. [35].

Alwassel *et al*, in their 2021 paper titled *TSP: Temporally Sensitive Pretraining of Video Encoders for Localization Tasks* introduce a supervised pre-training paradigm for TAL. This paradigm also considers background clips (which are not as important as for TAC) and global video information, to improve temporal sensitivity.

B.6.2 Contributions and Findings

- TSP trains an encoder to explicitly discriminate between foreground and background clips in untrimmed videos
- Temporally-sensitive features from TSP improves performance for TAL, action proposal and DVC Tasks
- Consistent performance gains for multiple algorithms, architectures and datasets
- TAL performance is boosted on short action instances

B.6.3 Pre-training Methodology

The major goal of TSP is to incorporate temporal sensitivity in video encoders. Hence, this pre-training involves training encoders on the task of classifying foreground clips, and classifying whether a clip is inside or outside the action. The data used for TSP consists of untrimmed videos with temporal annotations. The video encoder is trained end-to-end, from raw video input. From an untrimmed video, X is sampled, which has dimensions $3 \times L \times W \times H$, where L is the number of frames and W and H are the frame dimensions. 3 is the number of channels, i.e., RGB. Since there is a natural imbalance in the annotations of foreground and background clips, X is sampled in a way that an equal number of clips from each class is chosen. Each X is annotated with two labels:

- y^c : action class label, *if it is from a foreground clip*
- y^r : binary temporal region label, i.e. whether whether the clip is from a foreground/action region ($y^r = 1$) or background/no-action region ($y^r = 0$) of the video

B.6.4 Local and Global Feature Encoding

The encoder E transforms a clip X to a local F -dimensional feature f . Global Video Feature (GVF) is the max-pooled feature from all local features. For classifying the temporal region (foreground or background), combine each local feature f is combined with GVF; hence GVF acts as a conditional vector.

The video encoders used are ResNet3D and R(2+1)D, while the datasets used are ActivityNet v1.3 and THUMOS-14 [59]. The classification component involves two heads:

- An $F \times C$ fully connected layer for action classification: $f \rightarrow y^c$ (logits vector)
- A $2F \times 2$ fully connected layer for temporal region label (background foreground): $[f, GVF] \rightarrow y^r$

Cross-entropy loss is used for both classification heads. If the clip is a foreground clip, then loss from both heads are taken into account, relatively weighted by α^r and α^c . If the clip is a background clip, only the temporal region classification loss is taken, weighted by α^r .

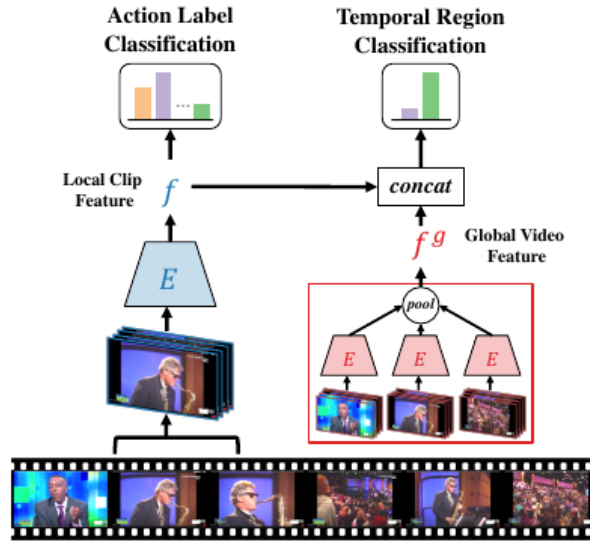


Figure B.9: The TSP framework. Image courtesy [35]

B.6.5 Performance

Using the pre-trained video encoder, Alwassel *et al* take some implementations of event localisation methods and feed temporally sensitive features to them. The algorithms used are:

- GTAD: sub-graph localization for Temporal Action Detection
- BMN: Boundary-matching network for Temporal Action Proposal Generation
- BMT: Bimodal Transformer for DVC [6] (audio features kept same)
- P-GCN: Graph CNN for Temporal Action Localization

Their observations are as follows [35]:

- Improved performance on multiple target tasks, such as TAC, TAL and DVC
- Consistency of performance regardless of type of video encoder used
- Consistent improvement in performance for all localization algorithms
- Indications of applicability of TSP on other datasets as well as its transferability

B.6.6 Conclusion

Alwassel *et al* present Temporally Sensitive Pretraining, a novel supervised pretraining approach for video encoders. This approach considers background clips along with foreground clips, and uses global information to gain temporal sensitivity. Their results indicate it is advantageous to use TSP features over other popular features to build accurate models [35], especially for DVC.

B.7 End-to-end Dense Video Captioning with Parallel Decoding

B.7.1 Overview

Wang *et al* in their paper, titled *End-to-end Dense Video Captioning with Parallel Decoding* [36] aim to address certain issues that they identify in existing DVC work, namely, the two-stage design and handcrafted and heuristic components for event proposals generation. They model DVC as a set prediction task, and propose a model which, using a transformer architecture based on the deformable DeTR [], performs event localization and event captioning in parallel. They utilize an event counter head to predict number of events based on video understanding.

B.7.2 Limitations Identified

Wang *et al* identify the following limitations of previous DVC work:

Two-stage design

Most models use a localize-then-describe strategy for DVC tasks, i.e. they first localize the events, and then feed these events to a captioning module for generating captions. These methods usually are not trained end-to-end; rather the two stages are trained separately. Consequently, the mutual promotion of the two subtasks is limited; captioning is considered a downstream task and hence inter-task associations are not learned.

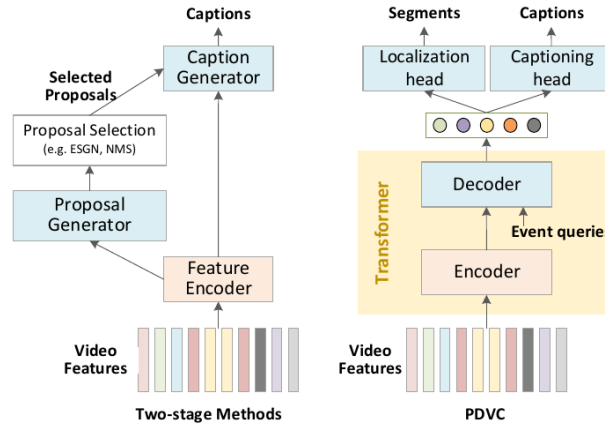


Figure B.10: Two-stage architecture versus proposed single-stage PDVC architecture. Image courtesy [36]

Unreliable Event Count Estimation

The number of event proposals is often heuristically determined, for example by applying manual thresholds on confidence scores, choosing the top k events, non-maximum suppression (NMS), et cetera. These methods introduce a lot of design issues, assumptions and hyperparameters, and the authors label these as *hand-crafted components* [36]. An unreliable event count estimation can cause either missing information in captions due to under-estimation, or redundancy in captions due to over-estimation.

Design of Proposal Generators

Most proposal generators use anchors and/or post-processing of events. Again, this brings in more hyperparameters and assumptions.

B.7.3 Proposed Solution

Wang *et al* model the task of dense video captioning as a *set prediction task* [36], i.e. the problem is to predict a set of $\{t_j^s, t_j^e, S_j\}$, where t_j^s and t_j^e are the start and end timestamps of the event j , and S_j is its caption.

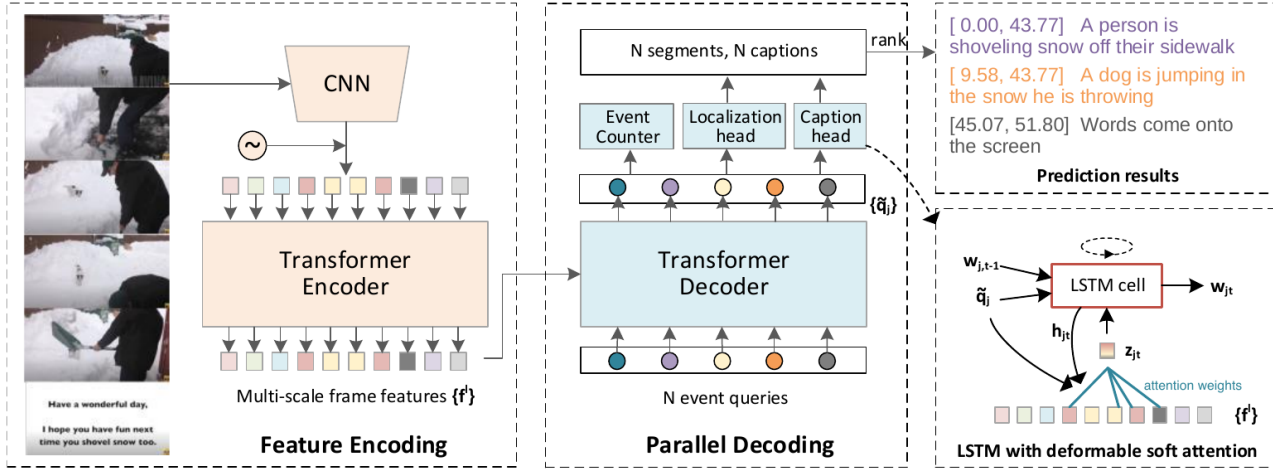


Figure B.11: PDVC Architecture. Image courtesy [36]

This set is of size N_{set} , which is predicted by the *event counter*. The event counter learns to predict the number of events in a video based on video understanding instead of heuristics.

The model performs the subtasks of event localization and caption generation in parallel, making the model a *single-stage pipeline* as opposed to the majority of previous work. This enables the end-to-end training of the entire model, and performance is boosted by mutual promotion of the two tasks.

The authors report localization performance being at par with state-of-the-art, while captioning quality is reported to be better than state-of-the-art at the time.

Feature Encoding

The features are extracted from videos using pretrained video encoders C3D and TSN. Temporal convolutional layers are applied to get *multi-scale features*, to get feature sequences across multiple resolutions [36]. These multi-scale features with positional encoding are given to the encoder of the deformable transformer, which uses *multi-scale deformable attention* to give a context vector as output.

Parallel Decoding

The decoder layers of the deformable transformer, which use multi-scale deformable attention in place of cross-attention (self-attention remaining unchanged as in [1]). The decoder layers query event-level features from the context vector conditioned on N learnable embeddings (termed *event queries*, q_j) and corresponding scalar reference point p_j . An event query q_j serves as an initial guess of event features and p_j serves as that of the center point of the event, and these are refined at each decoder layer. The output representation from decoder layers is given to the following three heads directly and in parallel.

Localization Head

This head performs box prediction with a center and length offset with respect to the reference point, and performs binary classification to output the foreground confidence for each event query. Both these predictions are performed using fully connected layers. The output of this head is $\{t_j^s, t_j^e, c_j^{loc}\}$, where c_j^{loc} is the localization (foreground) confidence.

Captioning Head

The authors propose two variants for the captioning head: (1) a vanilla LSTM and (2) an LSTM using deformable soft attention. The latter uses cues from the combination of caption words and the event features, while the former uses only event features.

Event Counter

This head consists of a max-pooling layer to compress event queries into a global feature vector, which is fed to a fully connected layer to predict vector r_{len} , from which N_{set} is found by the *argmax* function.

Finally, the top N_{set} events from N event queries in terms of their overall confidence scores c_j are selected.

$$c_j = c_j^{loc} + \mu \frac{1}{M_j^\gamma} \sum_{t=1}^{M_j} \log(c_{jt}^{cap})$$

where μ is a balancing factor between localization and captioning confidence, γ is a modulation factor to rectify the influence of caption length M_j .

Set Prediction Loss

The Hungarian algorithm is used to match predicted events with ground truths to find best bipartite matching results. The two sets are the predicted events and the ground truth events. Matching cost is defined as:

$$C = \alpha_{giou} L_{giou} + \alpha_{cls} L_{cls}$$

and the set prediction loss is defined as:

$$L = \beta_{giou} L_{giou} + \beta_{cls} L_{cls} + \beta_{ec} L_{ec} + \beta_{cap} L_{cap}$$

Here L_{giou} is the generalized IoU between predicted temporal segments and ground truth segments, L_{cls} is the focal loss between predicted classifications score and ground truth labels, L_{ce} is the cross-entropy between predicted event count distribution and ground truth and L_{cap} is the cross-entropy loss between predicted word distribution and ground truth (normalized by caption length).

B.8 End-to-End Object Detection with Transformers (DETR)

B.8.1 Overview

Carion *et al* in their 2020 paper titled *End-to-End Object Detection with Transformers* aims to streamline the detection pipeline by effectively removing the need for many hand-designed components by introducing a DEtection TRansformer or DETR which includes set-based global loss that forces unique predictions via bipartite matching, and a transformer encoder-decoder architecture.

B.8.2 DETR Model

Paper proposed two essential components for direct set predictions in detection:

- Set prediction loss forces unique matching between ground truth boxes and predicted boxes.
- An architecture that predicts a set of objects and models their relation in a single pass.

Object detection set prediction loss

- loss finds the optimal bipartite match between predicted and ground truth objects and then optimises object-specific (bounding box) losses.
- We pad the ground truth object with no objects if there are fewer objects in it to achieve one-to-one matching for direct set prediction without duplicates.

$$\hat{\sigma} = \arg \max_{\sigma \in \mathfrak{S}_N} \sum_i^N L_{match}(y_i, \hat{y}_{\sigma(i)})$$

- The next step is to calculate the Hungarian loss function for all pairings that were matched in the previous stage.

$$L_{Hungarian}(y, \hat{y}) = \sum_{i=0}^N [-\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbb{1}_{c_i \neq \phi} L_{box}(b_i, \hat{b}_{\hat{\sigma}(i)})]$$

DETR architecture

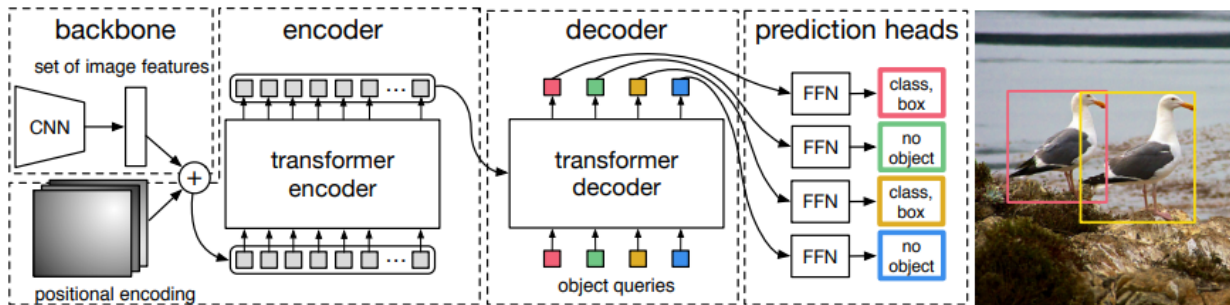


Figure B.12: DETR architecture by Carion *et al* (Courtesy [15])

Backbone

- It is a conventional CNN backbone which takes images as input and generates a lower-resolution activation map.

Transformer encoder

- As the encoder expects a sequence as input, spatial dimensions are collapsed into one dimension, resulting in a $d \times HW$ feature map.
- Each encoder layer has a standard architecture that includes a multi-head self-attention module and a feed forward network (FFN).
- Since the transformer architecture is permutation-invariant, we supplement it with fixed positional encodings that are added to the input of each attention layer.

Transformer decoder

- The decoder follows the transformer’s standard architecture, transforming N embeddings of size d with multi-headed self and encoder-decoder attention mechanisms.
- At each decoder layer, it decodes the N items in parallel. The N input embeddings must be distinct to produce different results because the decoder is permutation-invariant. These input embeddings are called object queries since they are learned positional encodings.
- The decoder converts the N object queries into an output embedding. A feed forward network decodes them independently into box coordinates and class labels, yielding N final predictions.

Prediction feed-forward networks

- The final prediction is computed by a 3-layer perceptron with ReLU activation function and hidden dimension d , and a linear projection layer.
- It predicts the box’s normalised centre coordinates, height, and width in relation to the input image, whereas the linear layer uses a softmax function to predict the class label.

B.8.3 Conclusion

Carion *et al* present DETR, a new design for object detection systems based on transformers and bipartite matching loss for direct set prediction. On the challenging COCO dataset, the technique produces results that are comparable to an optimised Faster R-CNN baseline. DETR is easy to set up and maintain, with a modular design that allows for easy expansion and competitive outcomes. Furthermore, it outperforms Faster R-CNN on huge objects.

B.9 Deformable DETR: Deformable Transformers for End-to-End Object Detection

B.9.1 Overview

Zhu *et al* in their 2021 paper titled *Deformable DETR: Deformable Transformers for End-to-End Object Detection* aims to replace the vanilla dense attention, which is the main computational bottleneck in DETR, with a deformable attention module. This can significantly reduce computational costs while also improving convergence.

B.9.2 Architecture

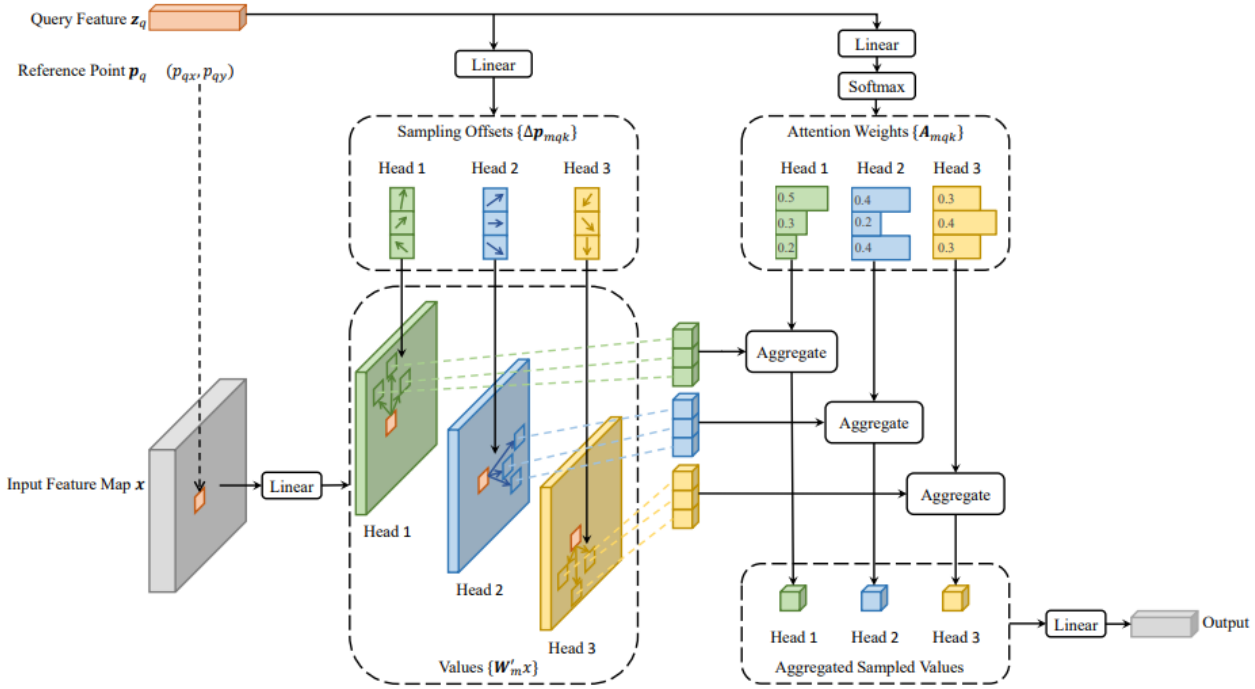


Figure B.13: Deformable Attention Module by Zhu *et al* (Courtesy [16])

Deformable Attention Module

- It attends to only a small set of key sampling points around a reference point, regardless of the spatial size of the feature maps. By assigning only a small fixed number of keys for each query, the issues of convergence and feature spatial resolution can be mitigated.
- Given an input feature map $x \in \mathbb{R}^{C \times H \times W}$, let q index a query element with content feature z_q and a 2-d reference point p_q , the deformable attention feature is calculated by

$$\text{DeformAttn}(z_q, p_q, x) = \sum_{m=1}^M W_m \left[\sum_{k=1}^K A_{mqk} \cdot W'_m x(p_q + \Delta p_{mqk}) \right]$$

where k indexes the sampled keys, m indexes the attention head, and K is the total sampled key number. Δp_{mqk} and A_{mqk} denote the sampling offset and attention weight of the k^{th} sampling point in the m^{th} attention head, respectively.

Multi-scale Deformable Attention Module

- Deformable attention module can be easily extended for multi-scale feature maps.
- Let $\{x^l\}_{l=1}^L$ be the input multi-scale feature maps, where $x^l \in \mathbb{R}^{C \times H^l \times W^l}$. Let $\hat{p}_q \in [0, 1]^2$ be the normalised coordinates of the reference point for each query element q , then the multi-scale deformable attention module is applied as

$$\text{DeformAttn}(z_q, \hat{p}_q, \{x^l\}_{l=1}^L) = \sum_{m=1}^M W_m \left[\sum_{l=1}^L \sum_{k=1}^K A_{mlqk} \cdot W'_m x^l(\phi_l(\hat{p}_q) + \Delta p_{mlqk}) \right]$$

where l indexes the input feature level, m indexes the attention head, and k indexes the sampling point. Δp_{mlqk} and A_{mlqk} denote the sampling offset and attention weight of the k th sampling point in the l th feature level and the m th attention head, respectively

Deformable Transformer Encoder

- The proposed multiscale deformable attention module replaces the transformer attention modules that process feature maps in DETR.
- The encoder's input and output are both multi-scale feature maps with the same resolutions.

Deformable Transformer Decoder

- The decoder includes cross-attention and self-attention modules.
- Object queries in the cross attention modules extract features from feature maps, where the key elements are from the encoder's output feature maps. Object queries interact with each other in the self-attention modules, where the key elements are of the object queries.
- Because the proposed deformable attention module is intended to process convolutional feature maps as key elements, the cross-attention module is replaced with the multi-scale deformable attention module, while the self-attention modules remain unaltered.

B.9.3 Conclusion

Zhu *et al* present Deformable DETR, an efficient and fast-converging end-to-end object detector. The multi-scale deformable attention modules, which are an efficient attention mechanism in processing image feature maps, are at the core of Deformable DETR.

B.10 Sparse DETR: Efficient End-to-End Object Detection with Learnable Sparsity

B.10.1 Overview

Roh *et al* in their 2021 paper *Sparse DETR: Efficient End-to-End Object Detection with Learnable Sparsity* aims to solve the bottleneck of Deformable DETR by sparsifying the encoder token by using DAM (Decoder Attention Map) predictor.

B.10.2 Methodology

Encoder Token Sparsification

In token sparsification scheme, the encoder module selectively refines a small number of encoder tokens. This encoder token subset is obtained from the backbone feature map x_{feat} with a certain criterion like Decoder Cross-Attention Map (DAM). For features that are not updated in this process, the values of x_{feat} are passed through the encoder layers without being changed

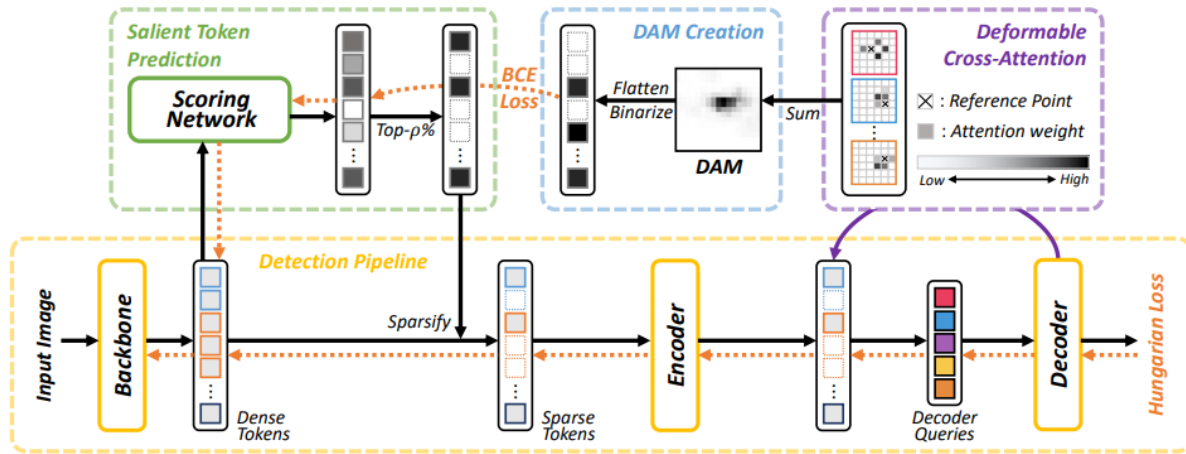


Figure B.14: Decoder Cross Attention Map (DAM) by Rho *et al* (Courtesy [17])

Decoder Cross-Attention Map

- A cross attention map from the transformer decoder is used to find the salient tokens.
- A scoring network is used to determine which encoder tokens should be further modified on the fly by predicting a pseudo groundtruth of the saliency defined by decoder cross-attention maps.
- The saliency of each input token is determined by aggregating the decoder cross-attentions between all object queries and the encoder output. It generates a single map of the same size as the backbone's feature map, which is known as the Decoder cross Attention Map (DAM).
- To train the scoring network, we binarize DAM so that only the top-k percentage of encoder tokens are kept. This is because the goal is to find a small subset of encoder tokens that the decoder references the most, rather than predicting how many times each encoder token will be referenced by the decoder.
- To predict how likely a given encoder token is to be included in the top-p % most referenced tokens, a 4-layer scoring network g is used, and the network is trained by minimising the binary cross entropy (BCE) loss between the binarized DAM and prediction.

$$L_{dam} = -\frac{1}{N} \sum_{i=1}^N BCE(g(X_{feat})_i, DAM_i^{bin})$$

Encoder Auxiliary Loss

Auxiliary loss can be applied to the encoder layers without sacrificing too much computational effort because the sparse detr model has a sparsified token set. Apart from better efficiency and performance, the encoder auxiliary loss has another advantage: it allows to stack more encoder layers safely without failing to converge.

B.10.3 Conclusion

Roh *et al* propose the encoder token sparsification method, which reduces the encoder’s attention complexity. They also propose a new sparsification criterion for selecting the informative subset from the complete token set which is Decoder Cross-Attention Map (DAM). Even when just 10% of the encoder token is used, Sparse DETR outperforms Deformable DETR and reduces overall computation by 38.

List of Figures

1.1	A Convolutional Network	7
1.2	The Transformer architecture. Image courtesy [1]	8
2.1	Evolution of dense video captioning methods over time.	12
3.1	Proposed model for Dense Video Captioning	21
3.2	ViViT encoder by Dosovitskiy <i>et al.</i> (Image courtesy [47])	22
3.3	AST encoder by Gong <i>et al.</i> (Image courtesy [48])	23
3.4	Distillation through attention by Touvron <i>et al.</i> (Image courtesy [50])	23
4.1	ViViT encoder by Dosovitskiy <i>et al.</i> (Image courtesy [47])	26
4.2	AST encoder by Gong <i>et al.</i> (Image courtesy [48])	27
4.3	The proposed feature extraction framework with ViViT and AST backbones. TSP heads for action classification and temporal region classification are shown on the right.	29
B.1	Pipeline introduced by Krishna <i>et al</i> (Image courtesy [2])	34
B.2	Dense Video Captioning Framework introduced by Iashin <i>et al</i> (Image courtesy [6])	35
B.3	Bi-modal Multiheaded Proposal Generator: Architecture (Image courtesy [6]	37
B.4	Video Transformer introduced by Dosovitskiy <i>et al</i> (Image courtesy [49])	39
B.5	Video Transformer introduced by Dosovitskiy <i>et al</i> (Courtesy [49])	39
B.6	Video Transformer introduced by Arnab <i>et al</i> (Courtesy [47])	42
B.7	Video Transformer introduced by Arnab <i>et al</i> (Image courtesy [47])	42
B.8	Audio Spectrogram Transformer by Gong <i>et al</i> (Courtesy [48])	44
B.9	The TSP framework. Image courtesy [35]	47
B.10	Two-stage architecture versus proposed single-stage PDVC architecture. Image courtesy [36]	48
B.11	PDVC Architecture. Image courtesy [36]	49
B.12	DETR architecture by Carion <i>et al</i> (Courtesy [15])	51
B.13	Deformable Attention Module by Zhu <i>et al</i> (Courtesy [16])	53
B.14	Decoder Cross Attention Map (DAM) by Rho <i>et al</i> (Courtesy [17])	55

List of Tables

2.1	Performance comparison of previous methods on the ActivityNet Captions dataset (* some videos unavailable)	14
2.2	Datasets and their characteristics	18

Bibliography

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017.
- [2] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *International Conference on Computer Vision (ICCV)*, 2017.
- [3] Nayyer Aafaq, Ajmal Mian, Wei Liu, Syed Zulqarnain Gilani, and Mubarak Shah. Video description. *ACM Computing Surveys*, 52(6):1–37, Jan 2020.
- [4] Yehao Li, Ting Yao, Yingwei Pan, Hongyang Chao, and Tao Mei. Jointly localizing and describing events for dense video captioning, 2018.
- [5] Jingwen Wang, Wenhao Jiang, Lin Ma, Wei Liu, and Yong Xu. Bidirectional attentive fusion with context gating for dense video captioning, 2018.
- [6] Vladimir Iashin and Esa Rahtu. A better use of audio-visual cues: Dense video captioning with bi-modal transformer, 2020.
- [7] Huijuan Xu, Boyang Li, Vasili Ramanishka, Leonid Sigal, and Kate Saenko. Joint event detection and description in continuous video streams, 2018.
- [8] Chaorui Deng, Shizhe Chen, Da Chen, Yuan He, and Qi Wu. Sketch, ground, and refine: Top-down dense video captioning. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 234–243, 2021.
- [9] Luowei Zhou, Yingbo Zhou, Jason J. Corso, Richard Socher, and Caiming Xiong. End-to-end dense video captioning with masked transformer. *arXiv preprint arXiv:1804.00819*, 2018.
- [10] Shyamal Buch, Victor Escorcia, Chuanqi Shen, Bernard Ghanem, and Juan Carlos Niebles. Sst: Single-stream temporal action proposals. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6373–6382, 2017.
- [11] Tianwei Lin, Xiao Liu, Xin Li, Errui Ding, and Shilei Wen. Bmn: Boundary-matching network for temporal action proposal generation, 2019.
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2020.
- [13] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows, 2021.
- [14] Mohsen Fayyaz, Soroush Abbasi Kouhpayegani, Farnoush Rezaei Jafari, Eric Sommerlade, Hamid Vaezi Joze, Hamed Pirsiavash, and Juergen Gall. Ats: Adaptive token sampling for efficient vision transformers. *arXiv:2111.15667*, December 2021.
- [15] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers, 2020.

- [16] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection, 2020.
- [17] Byungseok Roh, JaeWoong Shin, Wuhyun Shin, and Saehoon Kim. Sparse detr: Efficient end-to-end object detection with learnable sparsity, 2021.
- [18] Victor Escorcia, Fabian Caba Heilbron, Juan Carlos Niebles, and Bernard Ghanem. Daps: Deep action proposals for action understanding. In *ECCV*, 2016.
- [19] Luowei Zhou, Chenliang Xu, and Jason J. Corso. Towards automatic learning of procedures from web instructional videos, 2017.
- [20] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks, 2016.
- [21] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. Ssd: Single shot multibox detector. *Lecture Notes in Computer Science*, page 21–37, 2016.
- [22] Yilei Xiong, Bo Dai, and Dahua Lin. Move forward and tell: A progressive generator of video descriptions, 2018.
- [23] Yue Zhao, Yuanjun Xiong, Limin Wang, Zhirong Wu, Xiaoou Tang, and Dahua Lin. Temporal action detection with structured segment networks, 2017.
- [24] Huijuan Xu, Abir Das, and Kate Saenko. R-c3d: Region convolutional 3d network for temporal activity detection, 2017.
- [25] Xuguang Duan, Wenbing Huang, Chuang Gan, Jingdong Wang, Wenwu Zhu, and Junzhou Huang. Weakly supervised dense event captioning in videos, 2018.
- [26] Jonghwan Mun, Linjie Yang, Zhou Ren, Ning Xu, and Bohyung Han. Streamlined dense video captioning, 2019.
- [27] Tanzila Rahman, Bicheng Xu, and Leonid Sigal. Watch, listen and tell: Multi-modal weakly supervised dense event captioning, 2019.
- [28] Maitreya Suin and A. N. Rajagopalan. An efficient framework for dense video captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):12039–12046, Apr. 2020.
- [29] Vladimir Iashin and Esa Rahtu. Multi-modal dense video captioning, 2020.
- [30] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger, 2016.
- [31] Teng Wang, Huicheng Zheng, and Mingjing Yu. Dense-captioning events in videos: Sysu submission to activitynet challenge 2020, 2020.
- [32] Chuming Lin, Jian Li, Yabiao Wang, Ying Tai, Donghao Luo, Zhipeng Cui, Chengjie Wang, Jilin Li, Feiyue Huang, and Rongrong Ji. Fast learning of temporal action proposal via dense boundary generator, 2019.
- [33] Aman Chadha, Gurneet Arora, and Navpreet Kaloty. iperceive: Applying common-sense reasoning to multi-modal dense video captioning and video question answering, 2020.
- [34] Shaoxiang Chen and Yu-Gang Jiang. Towards bridging event captioner and sentence localizer for weakly supervised dense event captioning. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8421–8431, 2021.
- [35] Humam Alwassel, Silvio Giancola, and Bernard Ghanem. Tsp: Temporally-sensitive pretraining of video encoders for localization tasks, 2021.
- [36] Teng Wang, Ruimao Zhang, Zhichao Lu, Feng Zheng, Ran Cheng, and Ping Luo. End-to-end dense video captioning with parallel decoding, 2021.

- [37] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation, 2002.
- [38] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries, 2004.
- [39] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments, 2005.
- [40] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation, 2015.
- [41] Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. From word embeddings to document distances, 2015.
- [42] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation, 2016.
- [43] Hassan Akbari, Linagzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text, 2021.
- [44] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset, 2018.
- [45] Lorenzo Torresani Du Tran, Heng Wang and Matt Feiszli. Video classification with channel-separated convolutional networks, 2019.
- [46] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition, 2017.
- [47] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lucic, and Cordelia Schmid. Vivit: A video vision transformer, 2021.
- [48] Yuan Gong, Yu-An Chung, and James Glass. Ast: Audio spectrogram transformer. *arXiv preprint arXiv:2104.01778*, 2021.
- [49] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2020.
- [50] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.
- [51] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset, 2017.
- [52] Subhashini Venugopalan, Marcus Rohrbach, Jeff Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. Sequence to sequence – video to text, 2015.
- [53] Subhashini Venugopalan, Huijuan Xu, Jeff Donahue, Marcus Rohrbach, Raymond Mooney, and Kate Saenko. Translating videos to natural language using deep recurrent neural networks, 2015.
- [54] Haonan Yu, Jiang Wang, Zhiheng Huang, Yi Yang, and Wei Xu. Video paragraph captioning using hierarchical recurrent neural networks, 2016.
- [55] Shawn Hershey, Sourish Chaudhuri, Daniel P. W. Ellis, Jort F. Gemmeke, Aren Jansen, Channing Moore, Manoj Plakal, Devin Platt, Rif A. Saurous, Bryan Seybold, Malcolm Slaney, Ron Weiss, and Kevin Wilson. Cnn architectures for large-scale audio classification. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2017.

- [56] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
- [57] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Large scale learning of general visual representations for transfer. *CoRR*, abs/1912.11370, 2019.
- [58] Qizhe Xie, Eduard H. Hovy, Minh-Thang Luong, and Quoc V. Le. Self-training with noisy student improves imagenet classification. *CoRR*, abs/1911.04252, 2019.
- [59] Y.-G. Jiang, J. Liu, A. Roshan Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. THUMOS challenge: Action recognition with a large number of classes. <http://crcv.ucf.edu/THUMOS14/>, 2014.