Streamlined Dense Video Captioning

Jonghwan Mun^{1,5*} Linjie Yang² Zhou Ren³ Ning Xu⁴ Bohyung Han⁵ ¹POSTECH ²ByteDance AI Lab ³Wormpex AI Research ⁴Amazon Go ⁵Seoul National University

 1 jonghwan.mun@postech.ac.kr 2 linjie.yang@bytedance.com 3 zhou.ren@bianlifeng.com 4 ninxu@amazon.com 5 bhhan@snu.ac.kr

Abstract

Dense video captioning is an extremely challenging task since accurate and coherent description of events in a video requires holistic understanding of video contents as well as contextual reasoning of individual events. Most existing approaches handle this problem by first detecting event proposals from a video and then captioning on a subset of the proposals. As a result, the generated sentences are prone to be redundant or inconsistent since they fail to consider temporal dependency between events. To tackle this challenge, we propose a novel dense video captioning framework, which models temporal dependency across events in a video explicitly and leverages visual and linguistic context from prior events for coherent storytelling. This objective is achieved by 1) integrating an event sequence generation network to select a sequence of event proposals adaptively, and 2) feeding the sequence of event proposals to our sequential video captioning network, which is trained by reinforcement learning with two-level rewards-at both event and episode levels-for better context modeling. The proposed technique achieves outstanding performances on ActivityNet Captions dataset in most metrics.

1. Introduction

Understanding video contents is an important topic in computer vision. Through introduction of large-scale datasets [9, 31] and recent advances in deep learning technology, research towards video content understanding is no longer limited to activity classification or detection and addresses more complex tasks including video caption generation [1, 4, 13, 14, 15, 22, 23, 26, 28, 30, 33, 35, 36].

Video captions are effective for holistic video description. However, since videos usually contain multiple interdependent events in context of a video-level story (*i.e.* episode), a single sentence may not be sufficient to describe videos. Consequently, the dense video captioning



Episode: busking



Figure 1. An example of dense video captioning about a *busking* episode, which is composed of four interdependent events.

task [8] is introduced and getting more popular recently. Fig. 1 presents an example of dense video captioning for a *busking* episode, which is composed of four ordered events. Dense video captioning is conceptually more complex than simple video captioning since it requires detecting individual events in a video and understanding their context. Despite the complexity of the problem, most existing methods [8, 10, 27, 37] address the task as two sequential subtasks—event detection and event description—in which an event proposal network is in charge of detecting events and a captioning network generates captions for the selected proposals independently.

We propose a novel framework for dense video captioning, which considers the temporal dependency of the events. Contrary to existing approaches shown in Fig. 2(a), our algorithm detects event sequences from videos and generates captions sequentially, where each caption is conditioned on prior events and captions as illustrated in Fig. 2(b). Our algorithm has the following procedure. First, given a video, we obtain a set of candidate event proposals from an event proposal network. Then, an event sequence generation network selects a series of ordered events adaptively from the event proposal candidates. Finally, we generate captions for the selected event proposals using a sequential captioning network. The captioning network is trained via reinforcement learning using both event and episode-level rewards;

^{*}This work was done during the internship program at Snap Research.



Figure 2. Comparison between the existing approaches and ours for dense video captioning. Our algorithm generates captions for events sequentially conditioned on the prior ones by detecting an event sequence in a video.

the event-level reward allows to capture specific content in each event precisely while the episode-level reward drives all generated captions to make a coherent story.

The main contributions of the proposed approach are summarized as follows:

- We propose a novel framework detecting event sequences for dense video captioning. The proposed event sequence generation network allows the captioning network to model temporal dependency between events and generate a set of coherent captions to describe an episode in a video.
- We present reinforcement learning with two-level rewards, *episode* and *event* levels, which drives the captioning model to boost coherence across generated captions and the quality of description for each event.
- The proposed algorithm achieves state-of-the-art performances on the ActivityNet Captions dataset with large margins compared to the methods based on the existing framework.

The rest of the paper is organized as follows. We first discuss related works for our work in Section 2. The proposed method and its training scheme are described in Section 3 and 4 in detail, respectively. We present experimental results in Section 5, and conclude this paper in Section 6.

2. Related Work

2.1. Video Captioning

Recent video captioning techniques often adopt encoderdecoder frameworks inspired by success in image captioning [11, 16, 17, 25, 32]. Basic algorithms [22, 23] encode a video using Convolutional Neural Networks (CNNs) or Recurrent Neural Networks (RNNs), and convert the representation into a natural sentence using RNNs. Then various techniques are proposed to enhance the quality of generated captions by integrating temporal attention [33], joint embedding space of sentences and videos [14], hierarchical recurrent encoder [1, 13], attribute-augmented decoder [4, 15, 36], multimodal memory [28], and reconstruction loss [26]. Despite their impressive performances, they are limited to describing a video using a single sentence and can be applied only to a short video containing a single event. Thus, Yu *et al.* [35] propose a hierarchical recurrent neural network to generate a paragraph for a long video, while Xiong *et al.* [30] introduce a paragraph generation method based on event proposals, where an event selection module determines which proposals need to be utilized for caption generation in a progressive way. Contrary to these tasks, which generate a sentence or paragraph only for an input video, dense video captioning requires localizing and describing events at the same time.

2.2. Dense Video Captioning

Recent dense video captioning techniques [8, 10, 27, 37] are based on an identical framework, which attempts to solve the problem using two subtasks—event detection and caption generation; an event proposal network generates a set of candidate proposals and a captioning network is applied to each proposal. In this framework, the methods depend on a manual thresholding strategy to select final event proposals for description of input video contents.

Based on the framework, Krishna *et al.* [8] adopt a multiscale action proposal network [3], and introduce a captioning network that exploits visual context from past and future events with an attention mechanism. Wang *et al.* [27] employ a bidirectional RNN for improved event proposal generation, and propose a context gating mechanism in caption generation to adaptively control contribution of surrounding events. Li *et al.* [10] add temporal coordinate and descriptiveness regression for precise localization of event proposals, and adopt the attribute-augmented captioning network [34]. Rennie *et al.* [37] integrate a self-attention [20] for event proposal network and captioning network, and propose a masking network that converts the event proposals to differentiable masks and enables end-to-end learning of the two networks.

In contrast to the prior works, our algorithm identifies a small set of representative event proposals (*i.e.*, event sequences) for sequential caption generation, which enables us to generate coherent and comprehensive captions by exploiting both visual and linguistic context across selected events. However, the existing works only consider visual context, since the captioning network is applied to event proposals independently.

3. Our Framework

This section describes our main idea and the deep neural network architecture of our algorithm in detail.



Figure 3. Overall framework of the proposed algorithm. Given an input video, our algorithm first extracts a set of candidate event proposals $(p_1, p_2, p_3, p_4, p_5)$ using the Event Proposal Network (Section 3.2). From the candidate set, the Event Sequence Generation Network detects an event sequence $(\hat{e}_1 \rightarrow \hat{e}_2 \rightarrow \hat{e}_3)$ by selecting one out of the candidate event proposals (Section 3.3). Finally, the Sequential Captioning Network takes the detected event sequence and generates captions $(\hat{d}_1, \hat{d}_2, \hat{d}_3)$ conditioned on preceding events in a sequential manner (Section 3.4). The three models are trained in a supervised manner (Section 4.1) and then the Sequential Captioning Network is optimized additionally with reinforcement learning using two-level rewards (Section 4.2).

3.1. Overview

Let a video V contain a set of events $\mathcal{E} = \{e_1, \ldots, e_N\}$ with corresponding descriptions $\mathcal{D} = \{d_1, \ldots, d_N\}$, where N event are temporally localized using their starting and ending time steps. Existing algorithms [8, 10, 27, 37] typically divide the whole problem into two steps: event detection followed by description of detected events. These algorithms train models by minimizing sum of negative loglikelihoods of event and caption pairs as follows:

$$\mathcal{L} = \sum_{n=1}^{N} -\log p(d_n, e_n | V)$$
$$= \sum_{n=1}^{N} -\log p(e_n | V) p(d_n | e_n, V).$$
(1)

However, events in a video have temporal dependency and should be on a story about a single topic. Therefore, it is critical to identify an ordered list of events to describe a coherent story corresponding to the episode composed of the events. With this in consideration, we formulate dense video captioning as detection of an event sequence followed by sequential caption generation as follows:

$$\mathcal{L} = -\log p(\mathcal{E}, \mathcal{D}|V)$$

= $-\log p(\mathcal{E}|V) \prod_{n=1}^{N} p(d_n|d_1, \dots, d_{n-1}, \mathcal{E}, V).$ (2)

The overall framework of our proposed algorithm is illustrated in Fig. 3. For a given video, a set of candidate event proposals is generated by the Event Proposal Network. Then, our Event Sequence Generation Network provides a series of events by selecting one of candidate event proposals sequentially, where the selected proposals correspond to events comprising an episode in the video. Finally, we generate captions from the selected proposals using the proposed Sequential Captioning Network, where each caption is generated conditioned on preceding proposals and their captions. The captioning network is trained via reinforcement learning using event and episode-level rewards.

3.2. Event Proposal Network (EPN)

EPN plays a key role in selecting event candidates. We adopt single-stream temporal action proposals (SST) [2] due to its good performance and efficiency in finding semantically meaningful temporal regions via a single scan of videos. SST divides an input video into a set of nonoverlapping segments with a fixed length (i.e., 16 frames), where the representation of each segment is given by a 3D convolution (C3D) network [19]. By treating each segment as an ending point of an event proposal, SST identifies its matching starting points from the k preceding segments, which are represented by k-dimensional output vector from a Gated Recurrent Unit (GRU) at each time step. After extracting the top 1,000 event proposals, we obtain M candidate proposals, $\mathcal{P} = \{p_1, \ldots, p_M\}$, by eliminating highly overlapping ones using non-maximum suppression. Note that EPN provides representation of each proposal $p \in \mathcal{P}$, which is a concatenated vector of two hidden states at starting and ending segments in SST. This visual representation, denoted by Vis(p), is utilized for the other two networks.

3.3. Event Sequence Generation Network (ESGN)

Given a set of candidate event proposals, ESGN selects a series of events that are highly correlated and make up an episode for a video. To this ends, we employ a Pointer Network (PtrNet) [24] that is designed to produce a distribution over the input set using a recurrent neural network by adopting an attention module. PtrNet is well-suited for selecting an ordered subset of proposals and generating coherent captions with consideration of their temporal dependency.

As shown in Fig. 3, we first encode a set of candidate proposals, \mathcal{P} , by feeding proposals to an encoder RNN in an increasing order of their starting times, and initialize the first hidden state of PtrNet with the encoded representations to guide proposal selection. At each time step in PtrNet, we compute likelihoods a_t over the candidate event proposals and select a proposal with the highest likelihood out of all available proposals. The procedure is repeated until PtrNet happens to select the *END* event proposal, p_{end} , which is a special proposal to indicate the end of an event sequence.

The whole process is summarized as follows:

$$h_0^{\text{pur}} = \text{RNN}_{\text{enc}}(\text{Vis}(p_1), \dots, \text{Vis}(p_M)), \quad (3)$$

$$h_t^{\text{ptr}} = \text{RNN}_{\text{ptr}}(u(\hat{e}_{t-1}), h_{t-1}^{\text{ptr}}), \tag{4}$$

$$a_t = \operatorname{ATT}(h_t^{\operatorname{ptr}}, u(p_0), \dots, u(p_M)),$$
(5)

where h^{ptr} is a hidden state in PtrNet, ATT() is an attention function computing confidence scores over proposals, and the representation of proposal p in PtrNet, u(p) = [Loc(p); Vis(p)], is given by visual information Vis(p) as well as the location information Loc(p). Also, \hat{e}_t is a selected event proposal at time step t, which is given by

$$\hat{e}_t = p_{j^*}, \text{ where } j^* = \operatorname*{arg\,max}_{j \in \{0, \dots, M\}} a_t^j,$$
 (6)

where p_0 corresponds to p_{end} . Note that the location feature, Loc(p), is a binary mask vector, where the elements corresponding to temporal interval of the event are set to 1s and 0s otherwise. This is useful in identifying and disregarding proposals that overlap strongly with previously selected ones.

Our ESGN has clear benefits for dense video captioning. Specifically, it determines the number and order of events adaptively, which facilitates compact, comprehensive and context-aware caption generation. Noticeably, there are too many detected events in existing approaches (*e.g.*, \geq 50) given by manual thresholding. On the contrary, ESGN detects only 2.85 on average, which is comparable to the average number of events per video in ActivityNet Caption dataset, 3.65. Although sorting event proposals is an ill-defined problem, due to their two timestamps (starting and ending points), ESGN naturally learns the number and order of proposals based on semantics and contexts in individual videos by data-driven manner.

3.4. Sequential Captioning Network (SCN)

SCN employs a hierarchical recurrent neural network to generate coherent captions based on the detected event sequence $\hat{\mathcal{E}} = {\hat{e}_1, \ldots, \hat{e}_{N_s}}$, where $N_s (\leq M)$ is the number of selected events. As shown in Fig. 3, SCN is composed of two RNNs—an episode RNN and an event RNN, denoted by RNN ε and RNN $_e$, respectively. The episode RNN takes the proposals in the detected event sequence one by one and models the state of an episode implicitly, while the event RNN generates words in caption sequentially for each event proposal conditioned on the implicit representation of the episode, *i.e.*, based on the current context of the episode.

Formally, the caption generation process for the t^{th} event proposal in the detected event sequence is given by

$$r_t = \text{RNN}_{\mathcal{E}}(\text{Vis}(\hat{e}_t), g_{t-1}, r_{t-1}), \tag{7}$$

$$g_t = \text{RNN}_e(\text{C3D}(\hat{e}_t), \text{Vis}(\hat{e}_t), r_t), \quad (8)$$

where C3D(e) denotes feature descriptors of all segments within the span of event e based on C3D network, r_t is an episodic feature, and g_t means a feature (the last hidden state of RNN_e) of the generated caption from the t^{th} event proposal. The episode RNN provides the current episodic feature so that the event RNN generates context-aware captions, which are given back to the episode RNN.

Although both networks can be implemented with any RNNs conceptually, we adopt a single-layer Long Short-Term Memory (LSTM) with a 512 dimensional hidden state as the episode RNN, and a captioning network with temporal dynamic attention and context gating (TDA-CG) presented in [27] as the event RNN. TDA-CG generates words from a feature computed by gating a visual feature Vis(e) and an attended feature obtained from segment feature descriptors C3D(e).

Note that sequential captioning generation scheme enables to exploit both visual context (*i.e.* how other events look) and linguistic context (*i.e.* how other events are described) across events, and allows us to generate captions in an explicit context. Although existing methods [8, 27] also utilize context for caption generation, they are limited to visual context and model with no linguistic dependency due to their architectural constraints from independent caption generation scheme, which would result in inconsistent and redundant caption generation.

4. Training

We first learn the event proposal network and fix its parameters during training of other two networks. We train the event sequence generation network and the sequential captioning network in a supervised manner, and then further optimize the captioning network based on reinforcement learning with two-level rewards—at event and episode levels, respectively.

4.1. Supervised Learning

Event Proposal Network Let c_t^k be the confidence of the k^{th} event proposal at time step t in EPN—SST [2]. Denote the ground-truth label of the proposal by y_t^k , which is set to 1 if the event proposal has a temporal Intersection-over-Union (tIoU) with ground-truth events larger than 0.5, and 0 otherwise. Then, for a given video V and ground-truth labels y, we train EPN by minimizing a following weighted binary cross entropy loss:

$$\mathcal{L}_{\text{EPN}}(V, \mathcal{Y}) = -\sum_{t=1}^{T_c} \sum_{k=1}^{K} y_t^k \log c_t^k + (1 - y_t^k) \log(1 - c_t^k), \quad (9)$$

where $\mathcal{Y} = \{y_t^k | 1 \le t \le T_c, 1 \le k \le K\}$, K is the number of proposals containing each segment at the end and T_c is the number of segments in the video.

Event Sequence Generation Network For a video with ground-truth event sequence $\mathcal{E} = \{e_1, \ldots, e_N\}$ and a set of candidate event proposals $\mathcal{P} = \{p_1, \ldots, p_M\}$, the goal of ESGN is to select a proposal p highly overlapping with the ground-truth event e, which is achieved by minimizing the following sum of binary cross entropy loss:

$$\mathcal{L}_{\text{ESGN}}(V, \mathcal{P}, \mathcal{E}) = -\sum_{n=1}^{N} \sum_{m=1}^{M} \text{tIoU}(p_m, e_n) \log a_n^m \quad (10)$$
$$+ (1 - \text{tIoU}(p_m, e_n)) \log(1 - a_n^m),$$

where $IOU(\cdot, \cdot)$ is a temporal Intersection-over-Union value between two proposals, and a_n^m is the likelihood that the m^{th} event proposal is selected as the n^{th} event.

Sequential Captioning Network We utilize the groundtruth event sequence \mathcal{E} and its descriptions \mathcal{D} to learn our SCN via the *teacher forcing* technique [29]. Specifically, to learn two RNNs in SCN, we provide episode RNN and event RNN with ground-truth events and captions as their inputs, respectively. Then, the captioning network is trained by minimizing negative log-likelihood over words of ground-truth captions as follows:

$$\mathcal{L}_{cSeq}(V, \mathcal{E}, \mathcal{D}) = -\sum_{n=1}^{N} \log p(d_n | e_n)$$
(11)
= $-\sum_{n=1}^{N} \sum_{t=1}^{T_{d_n}} \log p(w_n^t | w_n^1, \dots, w_n^{t-1}, e_n),$

where $p(\cdot)$ denotes a predictive distribution over word vocabulary from the event RNN, and w_n^t and T_{d_n} mean the t^{th} ground-truth word and the length of ground-truth description for the n^{th} event.

4.2. Reinforcement Learning

Inspired by the success in image captioning task [16, 17], we further employ reinforcement learning to optimize the sequential captioning network. Similarly to the self-critical sequence training [17] approach, the objective of learning our captioning network is revised to minimize the negative expected rewards for sampled captions. The loss is formally given by

$$\mathcal{L}_{\text{SCN}}^{\text{RL}}(V, \hat{\mathcal{E}}, \hat{\mathcal{D}}) = -\sum_{n=1}^{N_s} \mathbb{E}_{\hat{d}_n} \left[\mathbf{R}(\hat{d}_n) \right], \qquad (12)$$

where $\hat{D} = \{\hat{d}_1, \dots, \hat{d}_{N_S}\}$ is a set of sampled descriptions from the detected event sequence $\hat{\mathcal{E}}$ with N_s events from ESGN, and $R(\hat{d})$ is a reward value for the individual sampled description \hat{d} . Then, the expected gradient on the sample set \hat{D} is given by

$$\nabla \mathcal{L}_{\text{SCN}}^{\text{RL}}(V, \hat{\mathcal{E}}, \hat{\mathcal{D}}) = -\sum_{n=1}^{N_s} \mathbb{E}_{\hat{d}_n} \left[\mathbf{R}(\hat{d}_n) \nabla \log p(\hat{d}_n) \right]$$
$$\approx -\sum_{n=1}^{N_s} \mathbf{R}(\hat{d}_n) \nabla \log p(\hat{d}_n). \tag{13}$$

We adopt a reward function with two levels: episode and event levels. This encourages models to generate coherent captions by reflecting the overall context of videos, while facilitating the choices of better word candidates in describing individual events depending on the context. Also, to reduce the variance of the gradient estimate [6, 16, 17], we use the rewards obtained from the captions generated with ground-truth proposals as baselines. This drives models to generate captions at least as competitive as the ones generated from ground-truth proposals, although intervals of event proposals are not exactly aligned with those of ground-truth proposals. Specifically, for a sampled event sequence $\hat{\mathcal{E}}$, we find a reference event sequence $\tilde{\mathcal{E}}$ = $\{\tilde{e}_1,\ldots,\tilde{e}_{N_0}\}$ and its descriptions $\tilde{\mathcal{D}} = \{\tilde{d}_1,\ldots,\tilde{d}_{N_0}\},\$ where the reference event \tilde{e} is determined to be one of ground-truth proposals with highest overlapping ratio with sampled event \hat{e} . Then, the reward for the n^{th} sampled description \hat{d}_n is obtained by

$$\mathbf{R}(\hat{d}_n) = (14) \\ \left[f(\hat{d}_n, \tilde{d}_n) - f(\check{d}_n, \tilde{d}_n) \right] + \left[f(\hat{\mathcal{D}}, \tilde{\mathcal{D}}) - f(\check{\mathcal{D}}, \tilde{\mathcal{D}}) \right],$$

where $f(\cdot, \cdot)$ returns a similarity score between two captions or two set of captions, and $\check{D} = \{\check{d}_1, \ldots, \check{d}_{N_s}\}$ denote the generated descriptions from the reference event sequence. Both terms in Eq. (14) encourage our model to increase the probability of sampled descriptions whose scores are higher than the results of generated captions from

Table 1. Event detection performances including recall and precision at four thresholds of temporal intersection of unions (@tIoU) on the ActivityNet Captions validation set. The bold-faced numbers mean the best performance for each metric.

Mathad		Re	ecall (@t	IoU)		Precision (@tIoU)				
wiethod	@0.3	@0.5	@0.7	@0.9	Average	@0.3	@0.5	@0.7	@0.9	Average
MFT [30]	46.18	29.76	15.54	5.77	24.31	86.34	68.79	38.30	12.19	51.41
ESGN (ours)	93.41	76.40	42.40	10.10	55.58	96.71	77.73	44.84	10.99	57.57

Table 2. Dense video captioning results including Bleu@N (B@N), CIDEr (C) and METEOR (M) for our model and other state-of-theart methods on ActivityNet Captions validation set. We report performances obtained from ground-truth (GT) proposals (left) and learned proposals (right). Asterisk (*) stands for methods re-evaluated on newer evaluation tool and star (\star) indicates methods exploiting additional modalities (*e.g.* optical flow, attribute) for video representation. The bold-faced numbers mean the best performance for each metric.

Mathad	with GT proposals					with learned proposals						
Wiethou	B@1	B@2	B@3	B@4	С	М	B@1	B@2	B@3	B@4	С	М
DCE [8]	18.13	8.43	4.09	1.60	25.12	8.88	10.81	4.57	1.90	0.71	12.43	5.69
DVC [10]*	19.57	9.90	4.55	1.62	25.24	10.33	12.22	5.72	2.27	0.73	12.61	6.93
Masked Transformer [37]**	23.93	12.16	5.76	2.71	47.71	11.16	9.96	4.81	2.42	1.15	9.25	4.98
TDA-CG [27]*	-	-	-	-	-	10.89	10.75	5.06	2.55	1.31	7.99	5.86
MFT [30]	-	-	-	-	-	-	13.31	6.13	2.82	1.24	21.00	7.08
SDVC (ours)	28.02	12.05	4.41	1.28	43.48	13.07	17.92	7.99	2.94	0.93	30.68	8.82

ground-truth event proposals. Note that the first and second term are computed on the current event and episode, respectively. We use two famous captioning metrics, METEOR and CIDEr, to define $f(\cdot, \cdot)$.

5. Experiments

5.1. Dataset

We evaluate the proposed algorithm on the ActivityNet Captions dataset [8], which contains 20k YouTube videos with an average length of 120 seconds. The dataset consists of 10,024, 4,926 and 5,044 videos for training, validation and test splits. The videos have 3.65 temporally localized events and descriptions on average, where the average length of the descriptions is 13.48 words.

5.2. Metrics

We use the performance evaluation tool¹ provided by the 2018 ActivityNet Captions Challenge, which measures the capability to localize and describe events². For evaluation, we measure recall and precision of event proposal detection, and METEOR, CIDEr and BLEU of dense video captioning. The scores of the metrics are summarized via their averages based on tIoU thresholds of 0.3, 0.5, 0.7 and 0.9 given identified proposals and generated captions. We use METEOR as the primary metric for comparison, since it is known to be most correlated to human judgments when only a small number of reference descriptions are available [21].

5.3. Implementation Details

For EPN, we use a two-layer GRU with 512 dimensional hidden states and generate 128 proposals at each ending segment, which makes the dimensionality of c_t in Eq. (9) 128. In our implementation, EPN based on SST takes a whole span of video for training as an input to the network, this allows the network to consider all ground-truth proposals, while the original SST [2] is trained with densely sampled clips given by the sliding window method.

For ESGN, we adopt a single-layer GRU and a singlelayer LSTM as EncoderRNN and RNN_{ptr} , respectively, where the dimensions of hidden states are both 512. We represent the location feature, denoted by $\text{Loc}(\cdot)$, of proposals with a 100 dimensional vector. When learning SGN with reinforcement learning, we sample 100 event sequences for each video and generate one caption for each event in the event sequence with a greedy decoding. In all experiments, we use Adam [7] to learn models with a mini-batch size of 1 video and a learning rate of 0.0005.

5.4. Comparison with Other Methods

We compare our Streamlined Dense Video Captioning (SDVC) algorithm with several existing state-of-theart methods including DCE [8], DVC [10], Masked Transformer [37] and TDA-CG [27]. We also report the results of MFT [30], which was originally proposed for video paragraph generation but its event selection module is also able to generate an event sequence from the candidate event proposals; it makes a choice between selecting a proposal for caption generation and skipping it, and constructs an event sequence implicitly. For MFT, we compare performances in both event detection and dense captioning.

Table 1 presents the event detection performances of

¹https://github.com/ranjaykrishna/densevid_eval ²On 11/02/2017, the official evaluation tool fixed a critical issue; only one out of multiple incorrect predictions for each video was counted. This leads to performance overestimation of [27, 37]. Thus, we received raw results from the authors and reported the scores measured by the new metric.

Table 3. Ablation results of me	ean averaged recall, precision	n and METEOR over four the	oU thresholds of 0.3, 0.5, 0.7	and 0.9 on the Activ-
ityNet Captions validation set.	We also present the number	of proposals in average. The	e bold-faced number means tl	he best performance.

Method	Proposal modules EPN ESGN		Captioning modules eventRNN episodeRNN		RL	Number of proposals	Recall	Precision	METEOR
EPN-Ind						77.99	84.97	28.10	4.58
ESGN-Ind		\checkmark	\checkmark			2.85	55.58	57.57	6.73
ESGN-SCN		\checkmark	\checkmark	\checkmark		2.85	55.58	57.57	6.92
ESGN-SCN-RL (SDVC)			\checkmark			2.85	55.58	57.57	8.82

Table 4. Results on ActivityNet Captions evaluation server.

	Audio	Flow	Visual	Ensemble	METEOR
RUC+CMU				yes	8.53
YH Technologies				no	8.13
Shandong Univ.			\checkmark	yes	8.11
SDVC (ours)				no	8.19

ESGN and MFT in ActivityNet Captions validation set. ESGN outperforms the progressive event selection module in MFT on most tIoUs with large margins, especially in recall. This validates the effectiveness of our proposed event sequence selection algorithm.

Table 2 illustrates performances of dense video captioning algorithms tested on ActivityNet Captions validation set. We measure scores with both ground-truth proposals and learned ones, where the number of predicted proposals in individual algorithms may be different; DCE, DVC, Masked Transformer and TDA-CG uses 1,000, 1,000, 226.78 and 97.61 proposals in average, respectively, while SDVC has only 2.85 proposals. According to Table 2, SDVC improves the quality of captions significantly compared to all other methods. Masked Transformer achieves comparable performance to ours using ground-truth proposals, but does not work well with learned proposals. Note that it uses optical flow features in addition to visual features, while SDVC is only trained on visual features. Since the motion information from optical flow features consistently improves the performances in other video understanding tasks [12, 18], incorporating motion information to our model may lead to performance improvement. MFT shows the highest METEOR score among existing methods, which is partly because MFT also considers temporal dependency across captions.

Table 4 presents the test split results from the evaluation server. SDVC achieves competitive performance based only on visual information while other methods exploit additional modalities (*e.g.*, audio and optical flow) to represent videos and/or perform model ensemble to boost accuracy as described in [5].

5.5. Ablation Studies

We perform several ablation studies on ActivityNet Captions validation set to investigate the contributions of individual components in our algorithm. In this experiment,

Table 5. Performance comparison varying reward levels in reinforcement learning on the ActivityNet Captions dataset.

Event-level reward	Episode-level reward	METEOR
		8.73
	\checkmark	8.29
\checkmark	\checkmark	8.82

we train the following four variants of our model: 1) *EPN-Ind*: generating captions independently from all candidate event proposals, which is a baseline similar to most existing frameworks, 2) *ESGN-Ind*: generating captions independently using eventRNN only from the events in the event sequence identified by our ESGN, 3) *ESGN-SCN*: generating captions sequentially using our hierarchical RNN from the detected event sequence, and 4) *ESGN-SCN-RL*: our full model (SDVC) which uses reinforcement learning to further optimize the captioning network.

Table 3 summarizes the results from this ablation study, and we have the following observations. First, the approach based on ESGN (ESGN-Ind) is more effective than the baseline that simply relies on all event proposals (EPN-Ind). Also, ESGN reduces the number of candidate proposals significantly, from 77.99 to 2.85 in average, with substantial increase in METEOR score, which indicates that ESGN successfully identifies event sequences from candidate event proposals. Second, context modeling through hierarchical structure (*i.e.*, event RNN + episode RNN) in a captioning network (ESGN-SCN) enhances performance compared to the method with independent caption generation without considering context (ESGN-Ind). Finally, large improvements in ESGN-SCN-RL indicates that reinforcement learning effectively improves the quality of captions for dense video captioning.

We also analyze the impact of two reward levels—event and episode—used for reinforcement learning. The results are presented in Table 5, which clearly demonstrates the effectiveness of training with rewards from both levels.

5.6. Qualitative Results

Fig. 4 illustrates qualitative results where the detected event sequences and generated captions are presented together. We compare the generated captions by our model (SDVC), which sequentially generate captions, with the model (ESGN-Ind) that generates descriptions indepen-



Figure 4. Qualitative results on ActivityNet Captions dataset. The arrows represent ground-truth events (red) and events in the predicted event sequence from our event sequence generation network (blue) for input videos. Note that the events in the event sequence are selected in the order of its index. For the predicted events, we show the captions generated independently (ESGN-Ind) and sequentially (SDVC). More consistent captions are obtained by our sequential captioning network, where words for comparison are marked in bold-faced black.

dently from the detected event sequences. The proposed ESGN successfully identifies event sequences for input videos. Also, sequential caption generation enables the model to describe events more coherently by exploiting both visual and linguistic contexts. For instance, in the first example on Fig. 4, SDVC can capture the linguistic context ('two men' in e_1 is represented by 'they' in both e_2 and e_3) as well as temporal dependency between events (an expression of 'continue' in e_3), while ESGN-Ind just recognizes and describes e_2 and e_3 as independently occurred events.

6. Conclusion

We presented a novel framework for dense video captioning, which considers visual and linguistic contexts for coherent caption generation by modeling temporal dependency across events in a video explicitly. Specifically, we introduced the event sequence generation network to detect a series of event proposals adaptively. Given the detected event sequence, a sequence of captions is generated by conditioning on preceding events in our sequential captioning network. We trained our models in a supervised manner while further optimizing the captioning network via reinforcement learning with two-level rewards for better context modeling. Our algorithm achieved the state-of-the-art performance in the ActivityNet Captions dataset on METEOR.

Acknowledgments This work was partly supported by Snap Inc., ASRI and Korean ICT R&D program of the MSIP/IITP grant [2016-0-00563, 2017-0-01780].

References

- Lorenzo Baraldi, Costantino Grana, and Rita Cucchiara. Hierarchical boundary-aware neural encoder for video captioning. In *CVPR*, 2017.
- [2] Shyamal Buch, Victor Escorcia, Chuanqi Shen, Bernard Ghanem, and Juan Carlos Niebles. Sst: Single-stream temporal action proposals. In *CVPR*, 2017.
- [3] Victor Escorcia, Fabian Caba Heilbron, Juan Carlos Niebles, and Bernard Ghanem. Daps: Deep action proposals for action understanding. In ECCV, 2016.
- [4] Zhe Gan, Chuang Gan, Xiaodong He, Yunchen Pu, Kenneth Tran, Jianfeng Gao, Lawrence Carin, and Li Deng. Semantic compositional networks for visual captioning. In *CVPR*, 2017.
- [5] Bernard Ghanem, Juan Carlos Niebles, Cees Snoek, Fabian Caba Heilbron, Humam Alwassel, Victor Escorcia, Ranjay Khrisna, Shyamal Buch, and Cuong Duc Dao. The activitynet large-scale activity recognition challenge 2018 summary. arXiv preprint arXiv:1808.03766, 2018.
- [6] Jiuxiang Gu, Jianfei Cai, Gang Wang, and Tsuhan Chen. Stack-captioning: Coarse-to-fine learning for image captioning. In AAAI, 2018.
- [7] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [8] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *ICCV*, 2017.
- [9] Yuncheng Li, Yale Song, Liangliang Cao, Joel Tetreault, Larry Goldberg, Alejandro Jaimes, and Jiebo Luo. TGIF: A New Dataset and Benchmark on Animated GIF Description. In *CVPR*, 2016.
- [10] Yehao Li, Ting Yao, Yingwei Pan, Hongyang Chao, and Tao Mei. Jointly localizing and describing events for dense video captioning. In *CVPR*, 2018.
- [11] Jonghwan Mun, Minsu Cho, and Bohyung Han. Textguided attention model for image captioning. In AAAI, 2017.
- [12] Phuc Nguyen, Ting Liu, Gautam Prasad, and Bohyung Han. Weakly supervised action localization by sparse temporal pooling network. In *CVPR*, 2018.
- [13] Pingbo Pan, Zhongwen Xu, Yi Yang, Fei Wu, and Yueting Zhuang. Hierarchical recurrent neural encoder for video representation with application to captioning. In *CVPR*, 2016.
- [14] Yingwei Pan, Tao Mei, Ting Yao, Houqiang Li, and Yong Rui. Jointly modeling embedding and translation to bridge video and language. In CVPR, 2016.

- [15] Yingwei Pan, Ting Yao, Houqiang Li, and Tao Mei. Video captioning with transferred semantic attributes. In CVPR, 2017.
- [16] Zhou Ren, Xiaoyu Wang, Ning Zhang, Xutao Lv, and Li-Jia Li. Deep reinforcement learning-based image captioning with embedding reward. In *CVPR*, 2017.
- [17] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jarret Ross, and Vaibhava Goel. Self-critical sequence training for image captioning. In CVPR, 2017.
- [18] Karen Simonyan and Andrew Zisserman. Twostream convolutional networks for action recognition in videos. In *NIPS*, 2014.
- [19] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015.
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- [21] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In CVPR, 2015.
- [22] Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. Sequence to sequence-video to text. In *ICCV*, 2015.
- [23] Subhashini Venugopalan, Huijuan Xu, Jeff Donahue, Marcus Rohrbach, Raymond Mooney, and Kate Saenko. Translating videos to natural language using deep recurrent neural networks. In NAACL-HLT, 2015.
- [24] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *NIPS*, 2015.
- [25] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In CVPR, 2015.
- [26] Bairui Wang, Lin Ma, Wei Zhang, and Wei Liu. Reconstruction network for video captioning. In CVPR, 2018.
- [27] Jingwen Wang, Wenhao Jiang, Lin Ma, Wei Liu, and Yong Xu. Bidirectional attentive fusion with context gating for dense video captioning. In *CVPR*, 2018.
- [28] Junbo Wang, Wei Wang, Yan Huang, Liang Wang, and Tieniu Tan. M3: Multimodal memory modelling for video captioning. In CVPR, 2018.
- [29] Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989.

- [30] Yilei Xiong, Bo Dai, and Dahua Lin. Move forward and tell: A progressive generator of video descriptions. In *ECCV*, 2018.
- [31] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *CVPR*, 2016.
- [32] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.
- [33] Li Yao, Atousa Torabi, Kyunghyun Cho, Nicolas Ballas, Christopher Pal, Hugo Larochelle, and Aaron Courville. Describing videos by exploiting temporal structure. In *ICCV*, 2015.
- [34] Ting Yao, Yingwei Pan, Yehao Li, Zhaofan Qiu, and Tao Mei. Boosting image captioning with attributes. In *ICCV*, 2017.
- [35] Haonan Yu, Jiang Wang, Zhiheng Huang, Yi Yang, and Wei Xu. Video paragraph captioning using hierarchical recurrent neural networks. In *CVPR*, 2016.
- [36] Youngjae Yu, Hyungjin Ko, Jongwook Choi, and Gunhee Kim. End-to-end concept word detection for video captioning, retrieval, and question answering. In *CVPR*, 2017.
- [37] Luowei Zhou, Yingbo Zhou, Jason J Corso, Richard Socher, and Caiming Xiong. End-to-end dense video captioning with masked transformer. In *CVPR*, 2018.