

# End-to-End Dense Video Captioning with Masked Transformer

Luowei Zhou\*  
University of Michigan  
luozhou@umich.edu

Yingbo Zhou\*  
Salesforce Research  
yingbo.zhou@salesforce.com

Jason J. Corso  
University of Michigan  
jjcorso@eecs.umich.edu

Richard Socher  
Salesforce Research  
richard@socher.org

Caiming Xiong<sup>†</sup>  
Salesforce Research  
cxiong@salesforce.com

# Motivation

- Dense Video Caption: event detection + event description.
- Previous models: build two models and train separately.
- Language information cannot have direct impacts on event proposal.
- This paper: end-to-end training, use language to help localization.
  - A **masking network** to convert discrete event proposals to differentiable mask, ensuring the consistency between proposal and caption.
  - **Self-attention** mechanism for capturing long-term dependencies.

# Method

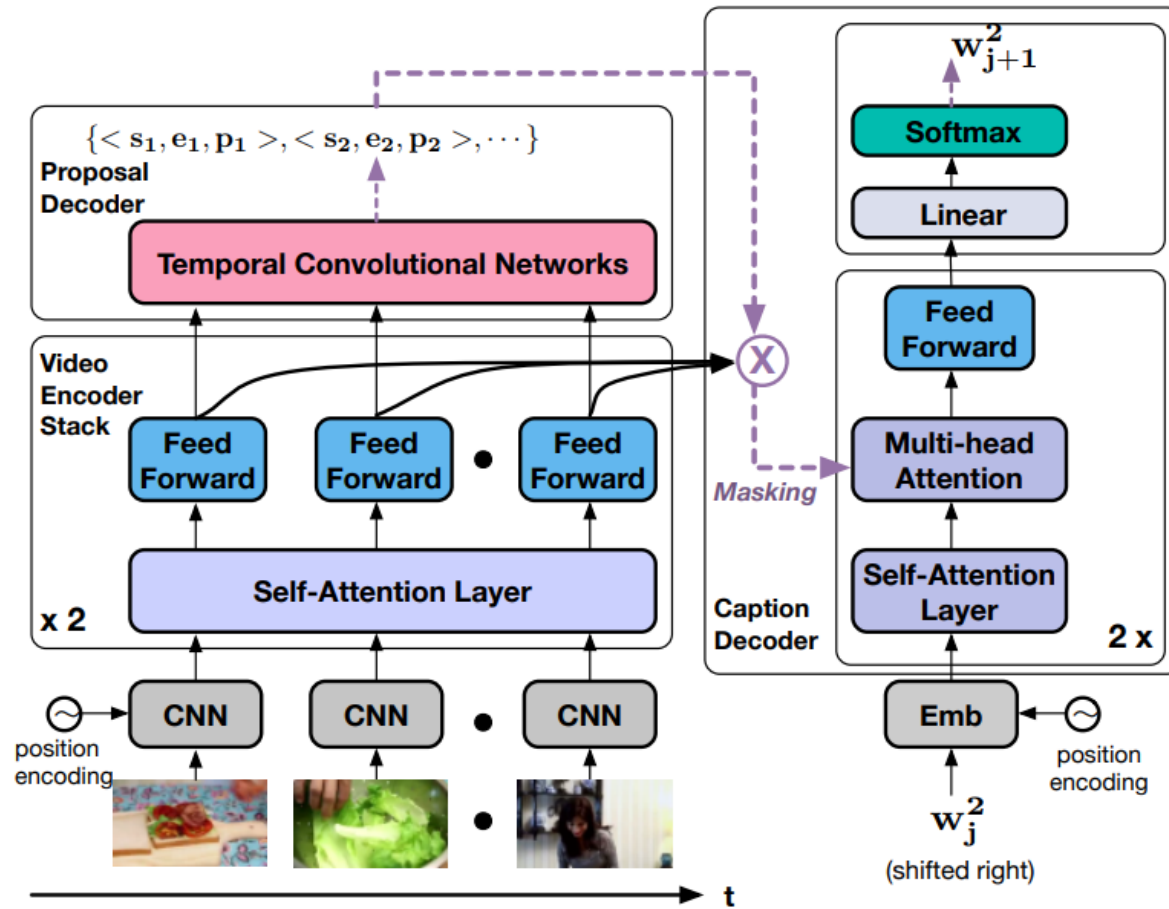


Figure 1. Dense video captioning is to localize (temporal) events from a video, which are then described with natural language sentences. We leverage temporal convolutional networks and self-attention mechanisms for precise event proposal generation and captioning.

# Method

- Video Encoder: a video  $X = \{x_1, \dots, x_T\} \rightarrow F$

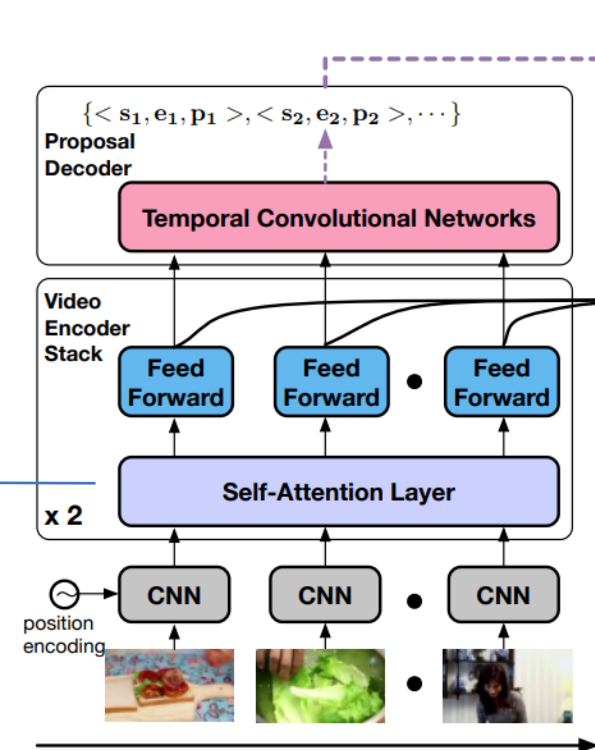
$$V(F^l) = \Psi(\text{PF}(\Gamma(F^l)), \Gamma(F^l))$$

$$\Gamma(F^l) = \begin{pmatrix} \Psi(\text{MA}(f_1^l, F^l, F^l), f_1^l)^\top & \dots \\ \Psi(\text{MA}(f_T^l, F^l, F^l), f_T^l)^\top \end{pmatrix}^\top$$

$$\Psi(\alpha, \beta) = \text{LayerNorm}(\alpha + \beta)$$

$$\text{PF}(\gamma) = M_2^l \max(0, M_1^l \gamma + b_1^l) + b_2^l$$

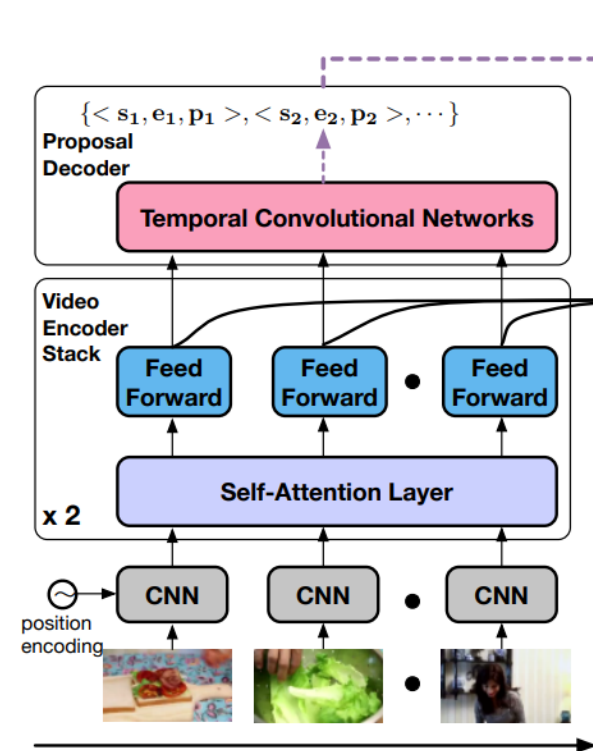
where  $\Psi(\cdot)$  represents the function that performs layer normalization on the residual output,  $\text{PF}(\cdot)$  denotes the 2-layered feed-forward neural network with ReLU nonlinearity for the first layer,  $M_1^l, M_2^l$  are the weights for the feed-



# Method

- Proposal Decoder: **N** explicit anchors.
  - Based on ProcNets\*, mainly temporal conv nets.
  - Each proposal is composed of score  $P_e$  and boundaries ( $S_p, E_p$ )

$$\begin{aligned} c_p &= c_a + \theta_c l_a & l_p &= l_a \exp\{\theta_l\}, \\ S_p &= c_p - l_p/2 & E_p &= c_p + l_p/2. \end{aligned}$$



\*: Towards Automatic Learning of Procedures from Web Instructional Videos, in AAAI 2018.

# Method

- Captioning Decoder: masked transformer.

$$Y_{\leq t}^{l+1} = C(Y_{\leq t}^l) = \Psi(\text{PF}(\Phi(Y_{\leq t}^l)), \Phi(Y_{\leq t}^l))$$

$$\Phi(Y_{\leq t}^l) = \begin{pmatrix} \Psi(\text{MA}(\Omega(Y_{\leq t}^l)_1, \hat{F}^l, \hat{F}^l), \Omega(Y_{\leq t}^l)_1)) \\ \dots \\ \Psi(\text{MA}(\Omega(Y_{\leq t}^l)_t, \hat{F}^l, \hat{F}^l), \Omega(Y_{\leq t}^l)_t)) \end{pmatrix}$$

Enc-Dec Attention

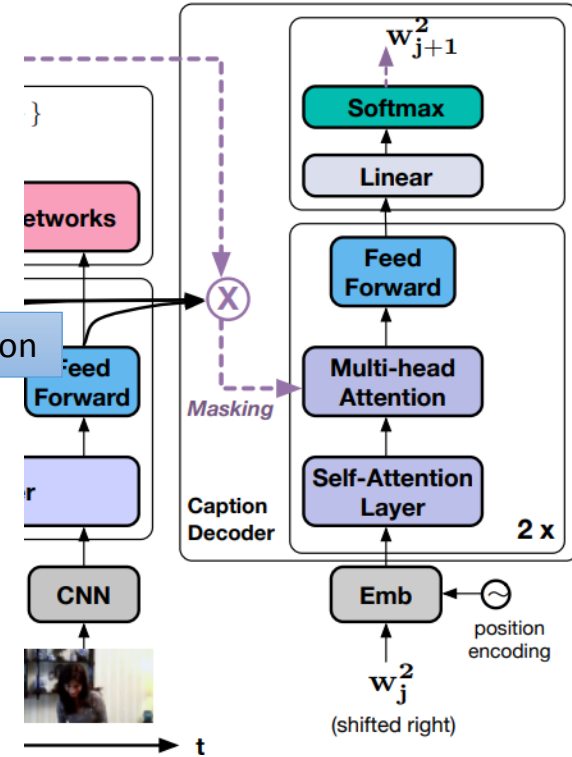
$$\Omega(Y_{\leq t}^l) = \begin{pmatrix} \Psi(\text{MA}(y_1^l, Y^l, Y^l), y_1^l)^\top \\ \dots \\ \Psi(\text{MA}(y_t^l, Y^l, Y^l), y_t^l)^\top \end{pmatrix}$$

Dec self-attention

$$\hat{F}^l = f_M(S_p, E_p) \odot F^l$$

Masking function

$$p(w_{t+1}|X, Y_{\leq t}^L) = \text{softmax}(W^V y_{t+1}^L)$$



# Method

- **Differentiable** Proposal Mask: proposal **specific** repre.

$$f_M(S_p, E_p, S_a, E_a, i) = \sigma(g([ \rho(S_p, :), \rho(E_p, :), \rho(S_a, :), \rho(E_e, :), \text{Bin}(S_a, E_a, :)])) \quad (1)$$

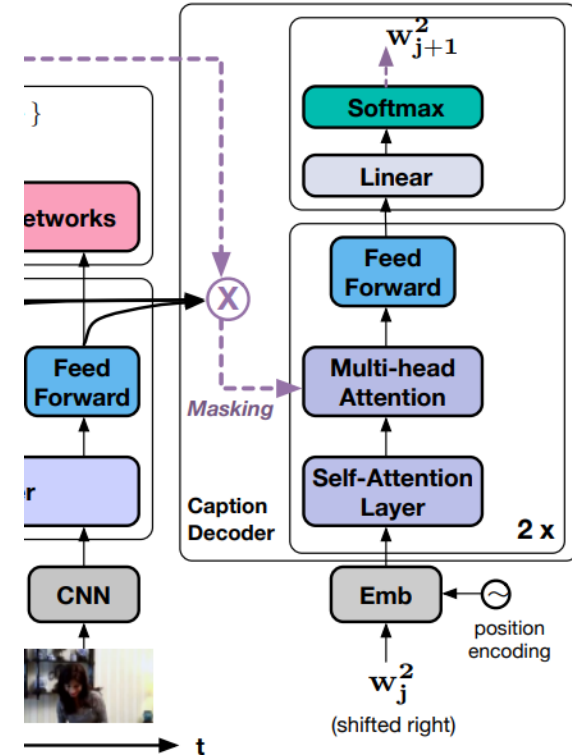
$$\rho(pos, i) = \begin{cases} \sin(pos/10000^{i/d}) & i \text{ is even} \\ \cos(pos/10000^{(i-1)/d}) & \text{otherwise} \end{cases} \quad (1)$$

$$\text{Bin}(S_a, E_a, i) = \begin{cases} 1 & \text{if } i \in [S_a, E_a] \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

- Gated masking:

$$f_{GM}(S_p, E_p, S_a, E_a, i) = P_e \text{Bin}(S_p, E_p, i) + (1 - P_e) f_M(S_p, E_p, S_a, E_a, i)$$

The continuous mask is used as a supplement for the proposal mask in case the **confidence** is low from the proposal module.



# Method

- Model training:

$$\mathcal{L}_r = \text{Smooth}_{\ell_1}(\hat{\theta}_c, \theta_c) + \text{Smooth}_{\ell_1}(\hat{\theta}_l, \theta_l)$$

Regression loss for event boundary.

$$\mathcal{L}_m^i = \text{BCE}(\text{Bin}(S_p, E_p, i), f_M(S_p, E_p, S_a, E_a, i))$$

Mask prediction loss.

$$\mathcal{L}_e = \text{BCE}(\hat{P}_e, P_e)$$

Event classification loss.

$$\mathcal{L}_c^t = \text{CE}(\hat{w}_t, p(w_t | X, Y_{\leq t-1}^L))$$

Captioning model loss.

$$\mathcal{L} = \lambda_1 \mathcal{L}_r + \lambda_2 \sum_i \mathcal{L}_m^i + \lambda_3 \mathcal{L}_e + \lambda_4 \sum_t \mathcal{L}_c^t$$



# Experiment

- Results from ActivityNet Caption Dataset:

Table 4. Event proposal results from ActivityNet Captions dataset. We compare our proposed methods with our baseline method ProcNets-prop on the validation set.

Method	Average Recall (%)
ProcNets-prop [42]	47.01
Bi-LSTM (ours)	50.65
Self-Attn (our)	<b>52.95</b>

Dense video captioning challenge leader board results. ts from the same team, we keep the highest one.

Method	METEOR
DEM [19]	4.82
Wang et al.	9.12
Jin et al.	9.62
Guo et al.	9.87
Yao et al. <sup>2</sup> (Ensemble)	12.84
Our Method	<b>10.12</b>

Table 1. Captioning results from ActivityNet Caption Dataset

Method	B@3	B@4	M
Bi-LSTM +TempoAttn	2.43	1.01	7.49
Masked Transformer	4.47	2.14	9.43
End-to-end Masked Transformer	<b>4.76</b>	<b>2.23</b>	<b>9.56</b>

# Summarization

- End-to-end training.
- First to use self-attention in video captioning.