

# PicSOM Team's Lessons Learned in TRECVID 2020 VTT

Jorma Laaksonen and Zixin Guo

Aalto University

Espoo, Finland

2020-12-08

# Overview

- Experiments with datasets and features
- Stacked attention captioning model
- Submissions and results
- Observations & issues

# Experiments with datasets and features

- We experimented with a selection training datasets ...
  - COCO – actually only images, but we used “fake” video features
  - TGIF
  - VATEX – new addition this year and we wanted to study its advantage
  - (MSR-VTT and MSVD had been used in earlier years, but dropped now)
- ... and features
  - ResNet-152
  - ResNext-101
  - I3D
  - C3D
  - (ResNet-101, semantic category, audio and multimodal features had been used earlier, but dropped now)

# Comparison of the datasets

<b>dataset</b>	<b>type</b>	<b>items</b>	<b>captions</b>
COCO	images	82783	414113
TGIF	videos	125713	125713
VATEX	videos	41250	825000

# “Fake 13D features” for COCO images

- 13D features can be extracted only from videos
- Average of 13D features of the T GIF videos were used as “fake 13D features” for COCO images
- Final input features were always concatenation of individual features
- Benefits of “fake 13D features” for the model training:
  - We can use also 414113 COCO captions
  - We can use genuine 13D features of T GIF and VATEX

# A selection of results on VTT 2019 ground truth data

CO CO	TG IF	VAT EX	Res Net	Res Next	I3D	C3D	METEOR	CIDEr	CIDErD	BLEU-4
X	X		X		X		0.2049	0.2348	0.1147	0.0319
X	X	X	X		X		0.2071	0.2746	0.1610	0.0443
X	X	X	X	X	X		0.2263	0.2812	0.1667	0.0446
	X	X	X	X	X		0.2253	0.2528	0.1518	0.0446
	X	X	X	X		X	0.2151	0.2345	0.1364	0.0397

# TRECVID 2019 result

CO CO	TG IF	VAT EX	Res Net	Res Next	I3D	C3D	METEOR	CIDER	CIDERD	BLEU-4
X	X		X		X		0.2049	0.2348	0.1147	0.0319
X	X	X	X		X		0.2071	0.2746	0.1610	0.0443
X	X	X	X	X	X		0.2263	0.2812	0.1667	0.0446
	X	X	X	X	X		0.2253	0.2528	0.1518	0.0446
	X	X	X	X		X	0.2151	0.2345	0.1364	0.0397

# Adding VATEX dataset to COCO and TGIF improves

CO CO	TG IF	VAT EX	Res Net	Res Next	I3D	C3D	METEOR	CIDEr	CIDErD	BLEU-4
X	X		X		X		0.2049	0.2348	0.1147	0.0319
X	X	X	X		X		<b>0.2071</b>	<b>0.2746</b>	<b>0.1610</b>	<b>0.0443</b>
X	X	X	X	X	X		0.2263	0.2812	0.1667	0.0446
	X	X	X	X	X		0.2253	0.2528	0.1518	0.0446
	X	X	X	X		X	0.2151	0.2345	0.1364	0.0397



# Adding ResNext features to ResNet and I3D improves

CO CO	TG IF	VAT EX	Res Net	Res Next	I3D	C3D	METEOR	CIDEr	CIDErD	BLEU-4
X	X		X		X		0.2049	0.2348	0.1147	0.0319
X	X	X	X		X		0.2071	0.2746	0.1610	0.0443
X	X	X	X	X	X		<b>0.2263</b>	<b>0.2812</b>	<b>0.1667</b>	<b>0.0446</b>
	X	X	X	X	X		0.2253	0.2528	0.1518	0.0446
	X	X	X	X		X	0.2151	0.2345	0.1364	0.0397

# Using COCO data with TGIF and VATEX improves

CO CO	TG IF	VAT EX	Res Net	Res Next	I3D	C3D	METEOR	CIDER	CIDERD	BLEU-4
X	X		X		X		0.2049	0.2348	0.1147	0.0319
X	X	X	X		X		0.2071	0.2746	0.1610	0.0443
X	X	X	X	X	X		<b>0.2263</b>	<b>0.2812</b>	<b>0.1667</b>	0.0446
	X	X	X	X	X		0.2253	0.2528	0.1518	0.0446
	X	X	X	X		X	0.2151	0.2345	0.1364	0.0397

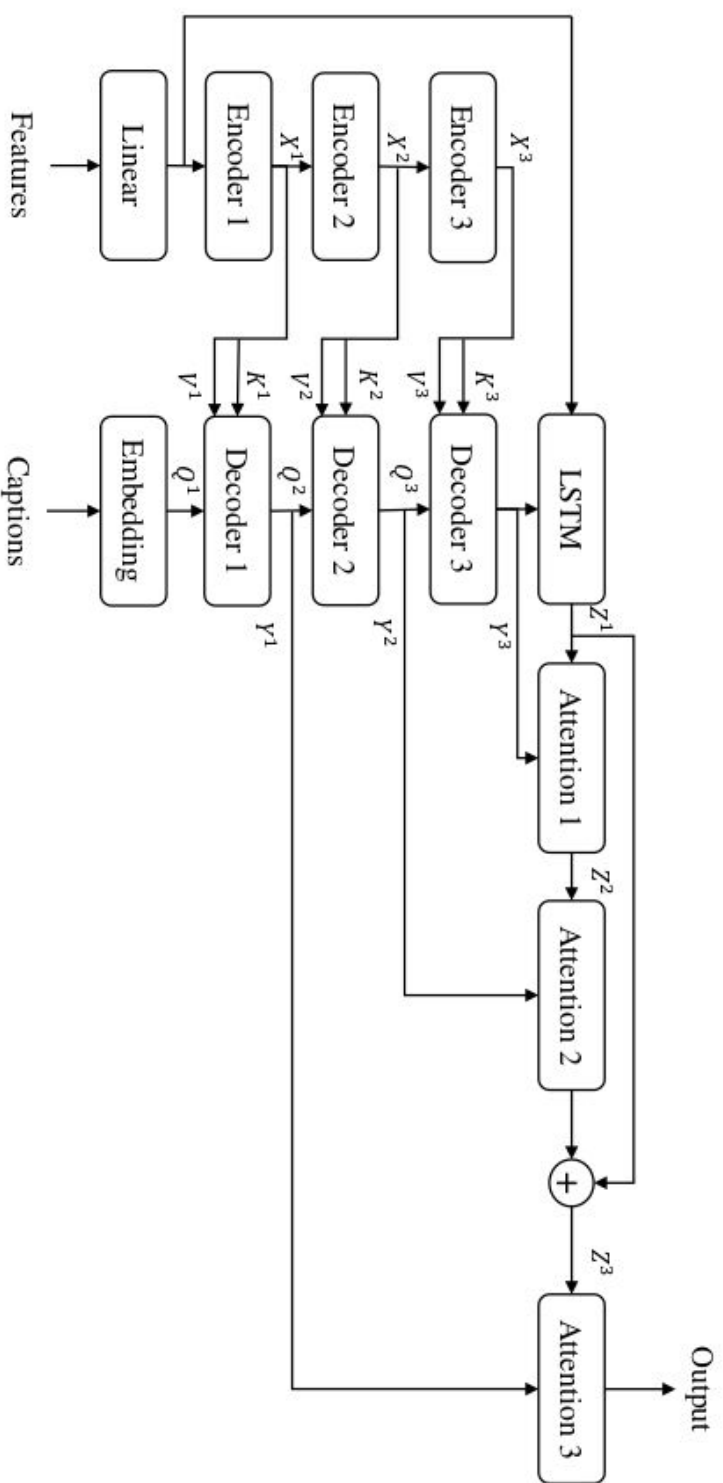
# I3D is better than C3D as a video feature

CO CO	TG IF	VAT EX	Res Net	Res Next	I3D	C3D	METEOR	CIDER	CIDERD	BLEU-4
X	X		X		X		0.2049	0.2348	0.1147	0.0319
X	X	X	X		X		0.2071	0.2746	0.1610	0.0443
X	X	X	X	X	X		0.2263	0.2812	0.1667	0.0446
	X	X	X	X	X		<b>0.2253</b>	<b>0.2528</b>	<b>0.1518</b>	<b>0.0446</b>
	X	X	X	X		X	0.2151	0.2345	0.1364	0.0397

# Final selection of datasets and features

CO CO	TG IF	VAT EX	Res Net	Res Next	I3D	C3D	METEOR	CIDEr	CIDErD	BLEU-4
X	X		X		X		0.2049	0.2348	0.1147	0.0319
X	X	X	X		X		0.2071	0.2746	0.1610	0.0443
X	X	X	X	X	X		0.2263	0.2812	0.1667	0.0446
	X	X	X	X	X		0.2253	0.2528	0.1518	0.0446
	X	X	X	X		X	0.2151	0.2345	0.1364	0.0397

# Stacked attention captioning model



# Stacked attention captioning model

- Based on the Transformer attention model

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_{\text{model}}}}\right)V$$

- Uses multihed attention

$$\text{Multihed}(Q, K, V) = \text{concat}(h_1, \dots, h_k)W^O$$

$$h_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \\ i = 1, \dots, k,$$

- Intra-modality stacked attention for visual features in the encoder
- Inter-modality stacked attention from visual to textual features in the decoder
- Stacked attention from decoder outputs to caption generation

$$\text{StackedAttention}(Y^{N-j+1}, Z^j) = \alpha(Y, Z) \odot Z \quad \alpha(Y, Z) = \sigma(W [Y, Z] + b)$$

# Submissions

run	description
1	Our latest and best <b>stacked attention</b> model, trained with <b>COCO+TGIF+VATEX</b>
2	Model similar to our best <b>VTT 2019 submission</b> , trained with <b>COCO+TGIF+VATEX</b>
3	Another well-performing <b>stacked attention</b> model, trained with <b>COCO+TGIF</b>
4	Model similar to our best <b>VTT 2019 submission</b> , trained with <b>COCO+TGIF</b>

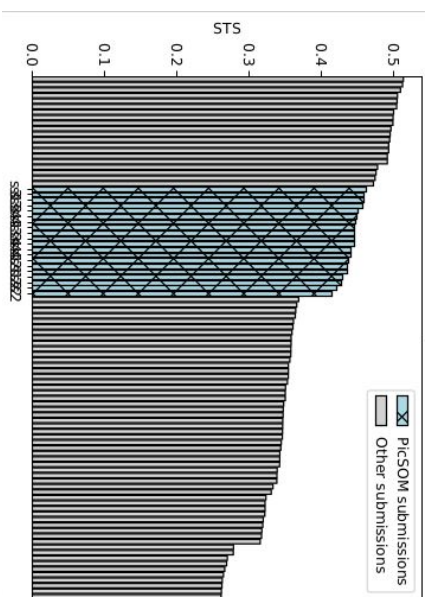
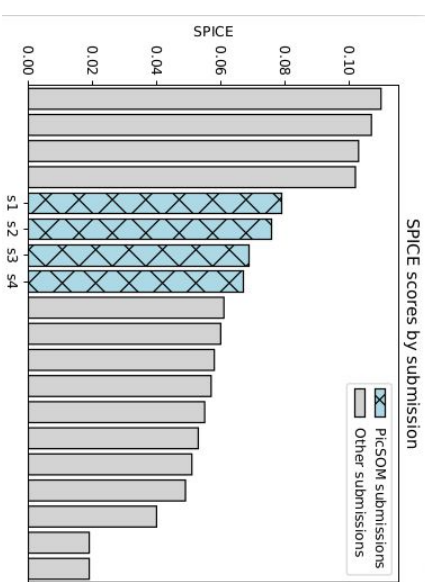
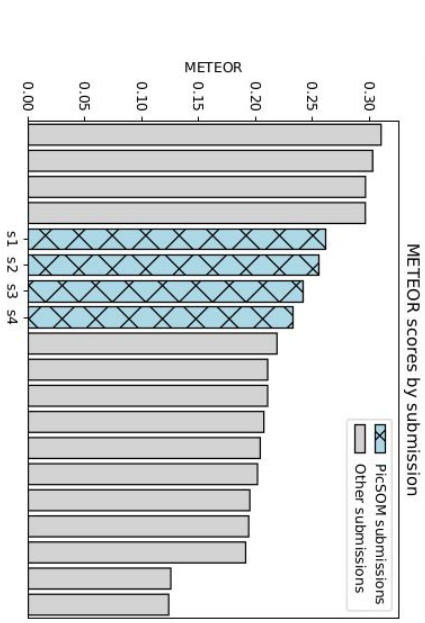
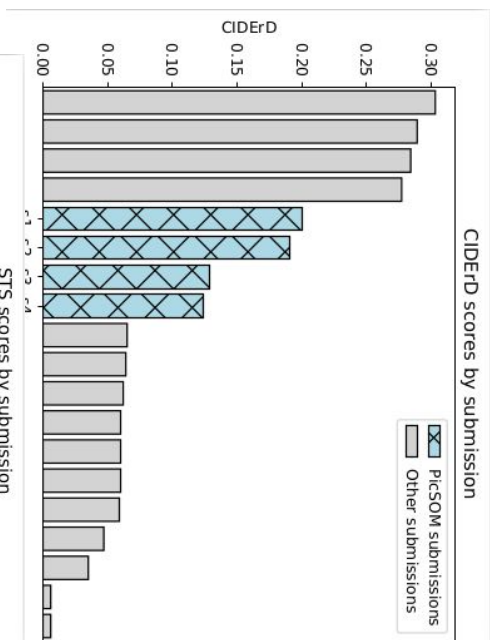
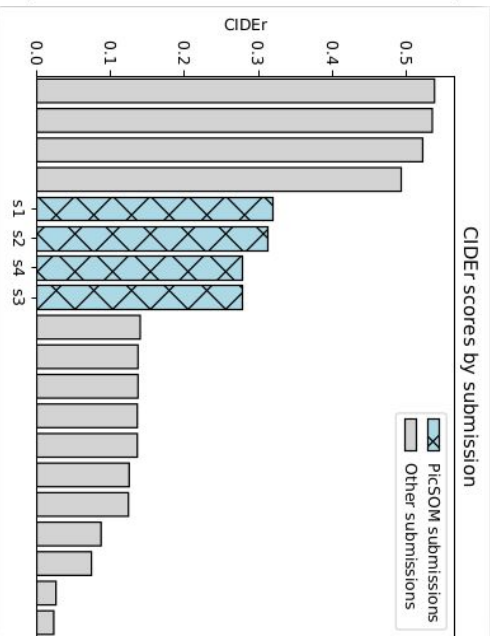
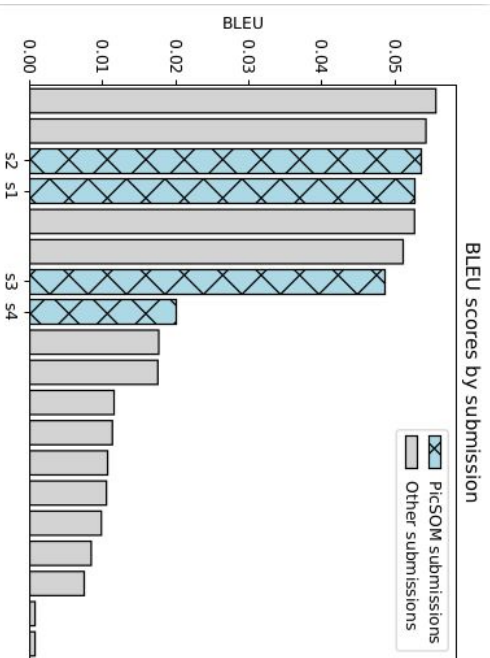
In all submissions cross-entropy training was used to initialize the LSTM model and self-critical reinforcement learning to finetune it with CLDER-D score. TRECVID VTT 2018 ground truth data were used for validation.

# Results

run	METEOR	CIDER	CIDERD	BLEU-4	SPICE	STS
1	<b>0.2617</b>	<b>0.319</b>	<b>0.200</b>	0.0527	<b>0.079</b>	0.4406
2	0.2556	0.312	0.191	<b>0.0536</b>	0.076	0.4293
3	0.2414	0.278	0.129	0.0485	0.069	<b>0.4581</b>
4	0.2323	0.278	0.124	0.0201	0.067	0.4458



# Plots



# Observations & issues

- Using VATEX data was more beneficial than the stacked attention model
- Current version of stacked attention didn't allow to use also ResNext features
- Our implementation of self-critical reinforcement learning is deficient as it creates captions in which the last word of the sentence is lacking
- Finding and correcting the bug will raise our scores substantially...