

## Motivation

### • Advantages of disentangled representations

- → Superior out-of-domain (OOD) generalization
- → Better interpretability
- → Better sample efficiency
- → Better transfer learning capabilities
- Use a diversity-enforcing loss to encourage disentangled representations.

### Background

- Tokenlearner
  - Adaptively learn a fixed set of token representations across one or more modalities.
  - Select a series of informative combinations of spatial locations in the image conditioned on all modalities.
  - For the  $i^{th}$  token  $z_i$ , it learns a spatial attention map  $\alpha_i(X)$  which is multiplied with the input X to generate a token output  $A_i(X)$ ,

$$z_i = A_i(X) = \rho(X \odot \gamma(\alpha_i(X)))$$

### Co-tokenization

- Cross-modality interaction during the visual feature extraction process by learning token representations, rather than rather considering them as an afterthought after feature extraction.
- Multiple streams of video at different spatio-temporal scales for multimodal representation learning.

### **Overall VideoQA Results**

Model	MSRVTT-QA	MSVD-QA	GFLOPs
Co-tokenization	33.7	32.5	67
Ours	33.1	30.1	41

**Table 1:** Comparison to state-of-the-art approaches for VideoQA (open 
 vocabulary). We pretrain on 10% subset of the HowTo69MVQA dataset, whereas Co-tokenization pretrained on the full HowTo100M dataset. We demonstrate competitive performance despite having a smaller model capacity.

# **Diversifying Joint Vision-Language Tokenization Learning**

\*Work done while at Google



# **Overall VQA Results**

Model	<b>GQA</b>	<b>SNLI-VE</b>
SimVLM (Huge)	_	86.32
UNITER	_	79.38
VinVL	65.05	
LXMERT	60.0	
Ours	76.79	80.15

**Table 2:** Comparison to state-of-the-art approaches (VQA)

Vardaan Pahuja<sup>1</sup>\*, AJ Piergiovanni<sup>2</sup>, Anelia Angelova<sup>2</sup> <sup>1</sup>The Ohio State University <sup>2</sup>Google DeepMind

# VideoQA Results

Dataset	<b>Pre-training</b>	Model	Accuracy
		Baseline	31.06
		Ours	31.37
MSRVTT-QA		Baseline	31.78
		Ours	33.05
		Baseline	27.98
		Ours	28.22
MSVD-QA		Baseline	28.08
		Ours	30.11
		Baseline	9.48
		Ours	9.96
IVQA		Baseline	8.86
		Ours	9.97

Table 3: Video QA results in the open vocabulary setting (val. set). The baseline is a similar capacity Co-tokenization model.

# VQA Results

		Val. set		Test set	
Dataset	Model	E.M.	<b>F1</b>	E.M.	<b>F1</b>
SNLI-VE	Baseline	76.70	76.70	76.59	76.59
	<b>Ours</b>	<b>78.06</b>	<b>78.06</b>	<b>77.36</b>	<b>77.37</b>
GQA	Baseline	73.48	73.56	<b>73.5</b>	<b>73.57</b>
	<b>Ours</b>	<b>75.02</b>	<b>75.11</b>	75.01	75.1

**Table 4:** VQA results in the pre-training setting.

		Val. set		Test set	
Dataset	Model	E.M.	<b>F1</b>	<b>E.M.</b>	<b>F1</b>
SNLI-VE	Baseline <b>Ours</b>	73.08 <b>73.15</b>	73.08 <b>73.15</b>	72.5 <b>72.69</b>	72.5 <b>72.69</b>
GQA	Baseline <b>Ours</b>	<b>68.08</b> 67.98	<b>68.13</b> 68.02	<b>68.14</b> 67.98	<b>68.2</b> 68.02

**Table 5:** VQA results in the no pre-training setting.







## Visualizations

• Localize attention to salient areas of the image, that are vital for answering the question.



### (a) Question Image

-		•		BK.	97. I		55
				BX.	SK.	BE I	
				ES ES	54	<b>83</b> 27.	BK.
				SI SI	85	20 20 20 20 20 20 20 20 20 20 20 20 20 2	20 20 20
	<b>(b)</b> T	oken	visual	ization	(Bas	eline	)
					2	55	75
							75
	( - )	<b>T</b> 1 .	•	-1:			

(C) IOKEN VISUALIZATION (OURS)

**Figure 1:** (a) *Question*: what kind of climbing vine or plant is this? Base*line*: tombppry, *Ours*: ivy, *Ground truth answers* = ['fern', 'grape', 'vine', 'ivy', 'unanswerable', 'creeping fig', 'unanswerable', 'unanswerable', 'ivy', 'green']; Bottom left: Weights assigned to each image patch for every token, lighter shades like yellow correspond to higher weights; *Bottom right*: Token attention masks grounded to the input image.