

# Unconditional Image-Text Pair Generation With Multimodal Cross Quantizer

Hyungyung Lee<sup>1</sup>, Sungjin Park<sup>1</sup>, Joonseok Lee<sup>2,3</sup>, Edward Choi<sup>1</sup>

<sup>1</sup>KAIST

<sup>2</sup>Google Research

<sup>3</sup>Seoul National University

Paper link:

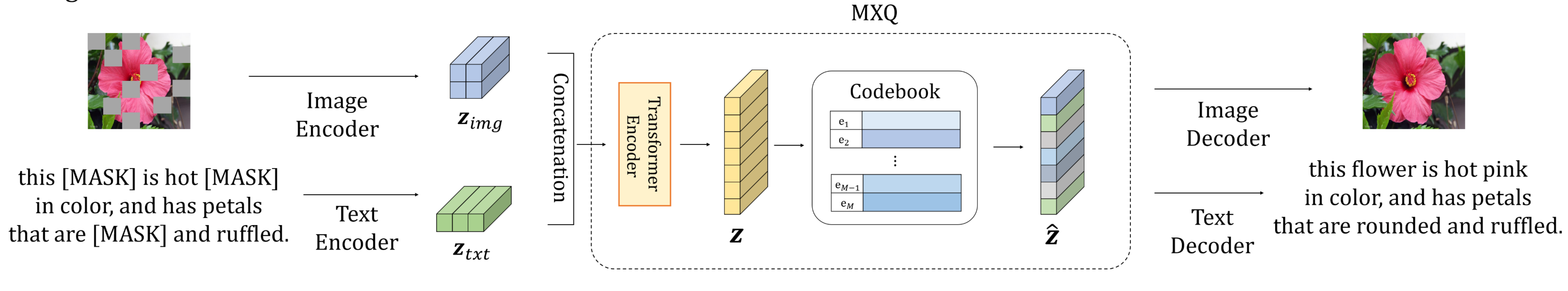


## Overview

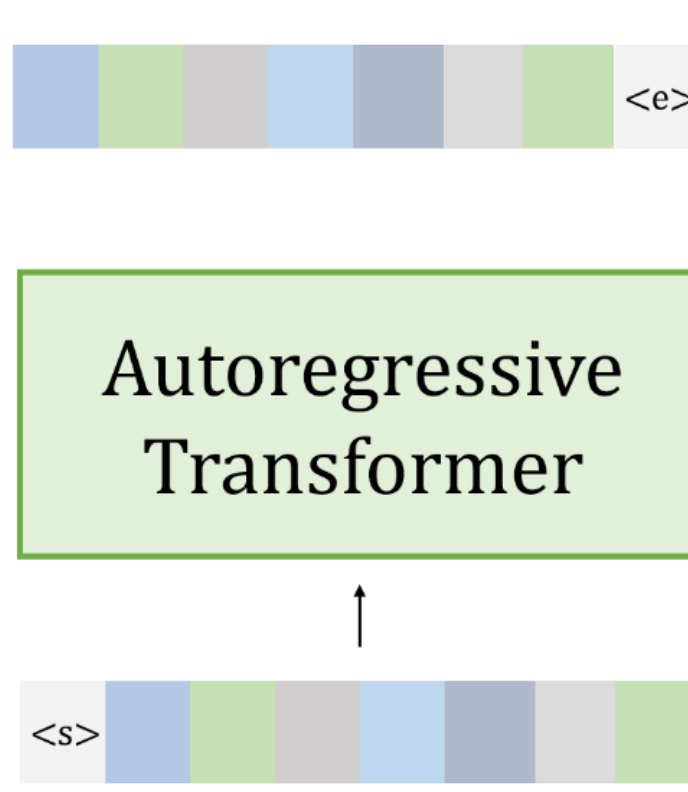
- Goal: we aim to generate image-text pairs simultaneously without any conditional input
- Contributions
  - We propose MXQ-VAE that learns a quantized joint representation space for unconditional image-text pair generation.
  - Experimental results reveal that MXQ-VAE generates a semantically consistent image-text pairs on multiple benchmark datasets
  - Also, MXQ-VAE learns meaningful semantic correlation between image and text in the quantized joint space.
  - We reveal that the quantized joint space leads to semantically consistent image-text pair generation.

## Proposed Method

### Stage 1



### Stage 2



- Stage 1: Learn a joint quantized representation space (named as MXQ-VAE)
  - MXQ-VAE takes masked image-text pairs as input, and learns a quantized joint representation space.
  - Then, the input is converted into a unified code sequence ( $\hat{z}$ ).
- Stage 2: Unconditional image-text pair generation with a unified discrete code sequence
  - Autoregressive transformer models the joint distribution over the code sequence.
  - At inference, MXQ-VAE decodes a sampled code sequence to an image-text pair.

## Dataset

- Caption MNIST (text augmented MNIST)
  - : Each image-text pair contains several color, digit and position.
- Oxford Flower-102, CUB-200-2011, COCO
- Degree Dataset
  - : To evaluate the semantic correlation between image and text in the quantized joint space in Stage 1.
  - : We gradually adjust the degree of the alignment between image and text by replacing the color and digit in text to other random color and digit.

Quadrant	Image	Text
Single		This is white 7.
Quad1		The 6 on the lower right is blue.

Quadrant	Image	Text
Quad2		The green 0 is on the upper right, and the top left 8 is red.
Quad3		the lower left 4 is blue, the red 0 is on the lower right, and the 2 on the top right is green.
Quad4		The bottom left 1 is blue, the upper right 7 is white, the bottom right 0 is red, and the upper left 3 is green.

Example of Caption MNIST

Image	Degree	Text
	2	The green 0 is on the upper right, and the top left 8 is red.
	1	The <u>white</u> 3 is on the upper right, and the top left 8 is red.
	0	The <u>blue</u> 1 is on the upper right, and the top left 7 is <u>green</u> .

Example of the Caption MNIST Quad2 Degree Dataset

Image	Degree	Text
	2	this flower has petals that are yellow and has brown stamen.
	1	this flower has petals that are <u>red</u> and has brown stamen.
	0	this flower has petals that are <u>maroon</u> and has <u>wine</u> stamen.

Example of the Flower Quad2 Degree Dataset

## Results

Dataset	Degree	Models	Degree 4	Degree 3	Degree 2	Degree 1	Degree 0
Caption MNIST	Quad4	Only Sharing $C$	0.486	0.443	0.394	0.358	0.315
		MXQ-VAE w/o TC	1.0	0.975	0.951	0.929	0.906
		MXQ-VAE w/o IM	0.896	0.802	0.698	0.595	0.498
		<b>MXQ-VAE (Ours)</b>	0.969	0.729	0.489	0.248	0.009
Flower	Quad4	Only Sharing $C$	0.939	0.704	0.516	0.321	0.131
		MXQ-VAE w/o TC	1.0	0.944	0.886	0.866	0.810
		MXQ-VAE w/o IM	0.997	0.728	0.482	0.278	0.067
		<b>MXQ-VAE (Ours)</b>	0.996	0.737	0.490	0.250	0.014
CUB	Quad4	Only Sharing $C$	0.985	0.771	0.572	0.356	0.155
		MXQ-VAE w/o TC	1.0	0.948	0.894	0.825	0.748
		MXQ-VAE w/o IM	0.998	0.833	0.645	0.424	0.181
		<b>MXQ-VAE (Ours)</b>	0.995	0.749	0.515	0.292	0.083

Table1. Multimodal semantic correlation on the Degree datasets

- Enhance the alignment between image and text (Table 1)

→ Measure reconstruction accuracy between the reconstructed text and the input text of the Degree dataset

→ Our model reaches near 1.0, 0.75, 0.5, 0.25, 0.0 on the Quad4 Degree dataset.

- Generate semantically consistent image-text pairs (Table 2, 3)

→ Measure semantic consistency between the generated image and text

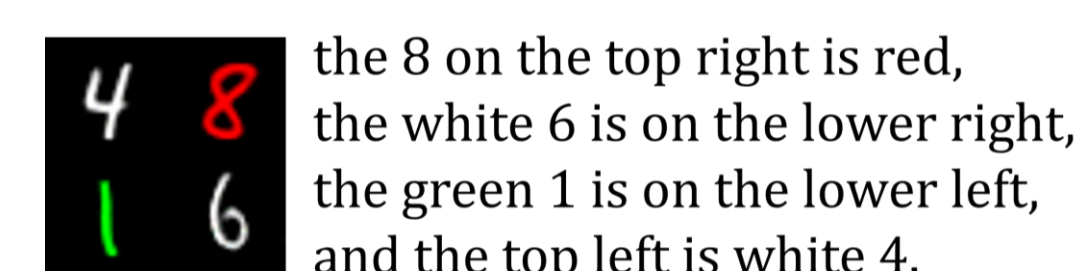
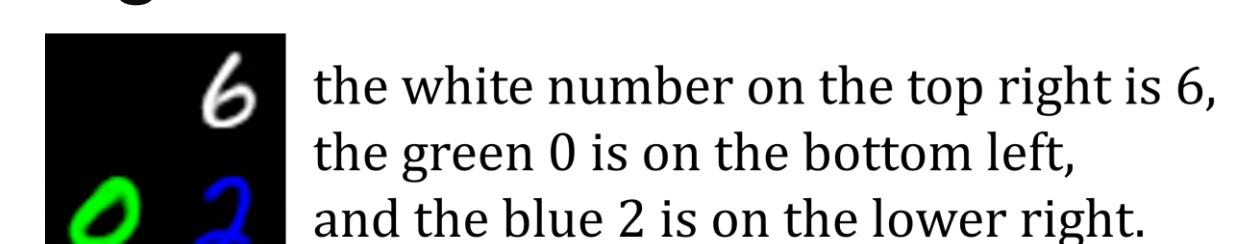
→ Our model outperforms all baselines on three datasets.

Models	Single	Quad1	Quad2	Quad3	Quad4	Average
I&T	0.979	0.926	0.675	0.434	0.255	0.654
T&I	0.803	0.780	0.458	0.282	0.161	0.497
$I_{T\_Emb}$	0.953	0.953	0.956	0.958	0.849	0.945
$T_{Emb\_I}$	0.086	0.895	0.913	0.916	0.828	0.728
<b>MXQ-VAE (Ours)</b>	<b>0.998</b>	<b>0.997</b>	<b>0.994</b>	<b>0.996</b>	<b>0.974</b>	<b>0.992</b>

Table2. Semantic consistency of the generated samples on Caption MNIST

Models	Flower				CUB			
	Modified unigram	Sentence similarity			Modified unigram	Sentence similarity		
	precision	Top-1	Top-5	Top-10	precision	Top-1	Top-5	Top-10
Joint GAN [10] + MMVAE [11]	0.324	0.808	0.788	0.774	-	-	-	-
<b>MXQ-VAE (Ours)</b>	<b>0.428</b>	<b>0.941</b>	<b>0.926</b>	<b>0.916</b>	<b>0.478</b>	<b>0.948</b>	<b>0.919</b>	<b>0.900</b>

Table3. Semantic consistency of the generated samples on Flower and CUB



Generated image-text pairs