# My Research Topic

related to natural language processing

2023-11-30

Hwang Seunghyeon | Augmented Intelligence Lab

# Contents

- What is image capturing?
  - LG Captioning AI
  - Video captioning
- Real-time video captioning with MR devices
  - Zero-Shot Image Captioning
  - Research and implementation goals

# Image Captioning

Image to Text

Input: photos or drawings

Output: sentences or keywords

CVPR 2023 - LG Captioning AI

# Video Captioning

Generate captions that match the words and actions spoken on the screen.

Designed to be used by people with hearing impairments.



**YouTube experiments with AI auto-generated video summaries**

Lauren Forristal

@laurenforristal / 11:56 PM GMT+9 • August 1, 2023

Comment

# Zero-shot Image Captioning

# ZeroCap

Yoad Tewel, Yoav Shalev,
Idan Schwartz, Lior Wolf

School of Computer Science,
Tel Aviv University

Using a visual-semantic model
with a large language model to
generate caption

visual-semantic model : CLIP

large language model : GPT-2

# ZeroCap: Zero-Shot Image-to-Text Generation for Visual-Semantic Arithmetic

*Yoad Tewel, Yoav Shalev, Idan Schwartz, Lior Wolf*, Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022, pp. 17918-17928

## Abstract

Recent text-to-image matching models apply contrastive learning to large corpora of uncurated pairs of images and sentences. While such models can provide a powerful score for matching and subsequent zero-shot tasks, they are not capable of generating caption given an image. In this work, we repurpose such models to generate a descriptive text given an image at inference time, without any further training or tuning step. This is done by combining the visual-semantic model with a large language model, benefiting from the knowledge in both web-scale models. The resulting captions are much less restrictive than those obtained by supervised captioning methods. Moreover, as a zero-shot learning method, it is extremely flexible and we demonstrate its ability to perform image arithmetic in which the inputs can be either images or text and the output is a sentence. This enables novel high-level vision capabilities such as comparing two images or solving visual analogy tests. Our code is available at: https://github.com/YoadTew/zero-shot-image-to-text.

# Zero-shot Image Captioning

- AI makes inferences about objects it sees <span style="color:yellow">for the first time.</span>
  - Based on pre-trained data.
  - Just like a human.
- Combining
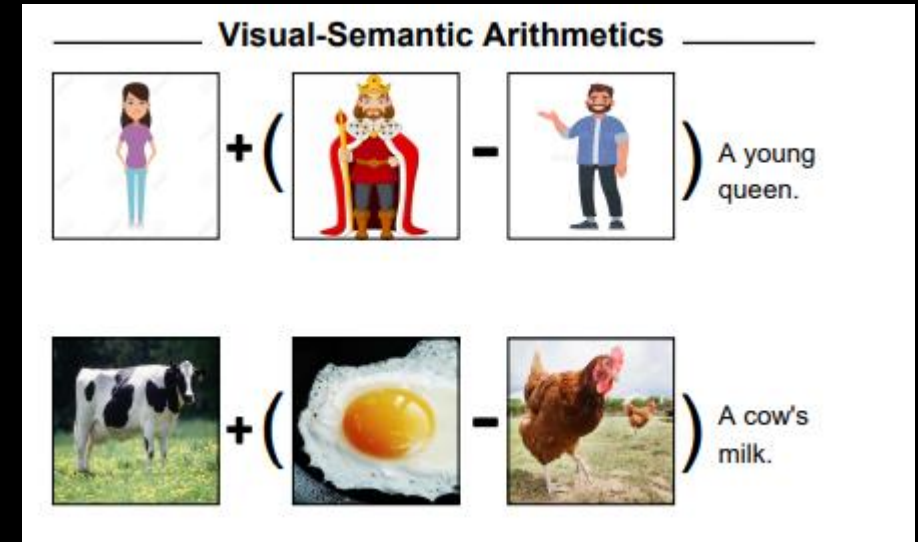  - visual-semantic model.
  - Large language model. (LLM)

# LLM

- Inferring the word that follows the caption

  - Image of a ……

- Inference Methods

  - Generate words in the correct orientation for a given image

    - CLIP loss $L_{CLIP}$

  - Preserves language attributes

    - loss term $L_{CE}$

# Visual-Semantic Arithmetic

- Word Embedding arithmetic
  - queen ='king'-'man'+ 'woman'
  - Embedding words as vectors
- Image Embedding arithmetic
  - Embedding Vectors in Image Space

# Comparison

## Traditional Captioning

- Pre-trained models

- Doesn't work well with new data

## Zero-shot Captioning

- LLM

- No need to learn specific visuals

# Real-time video captioning with MR devices

Apple Vision Pro applications for the deaf and hard of hearing

# Research goals

- Enabling real-time video captioning with zero-shot captioning technology

- Combining **multiple** state-of-the-art
  - Visual-semantic model using Diffusion Models
  - GPT-4

- Analyze usability differences with existing products

# Implementation goals

- Purpose
  - Real-time captioning in movie theaters, performance venues, etc.
  - Filling the gap that OTT cannot fill for the hearing impaired.
- Environment
  - Apple Vision Pro
  - Microsoft Holo Lens 2

# Existing

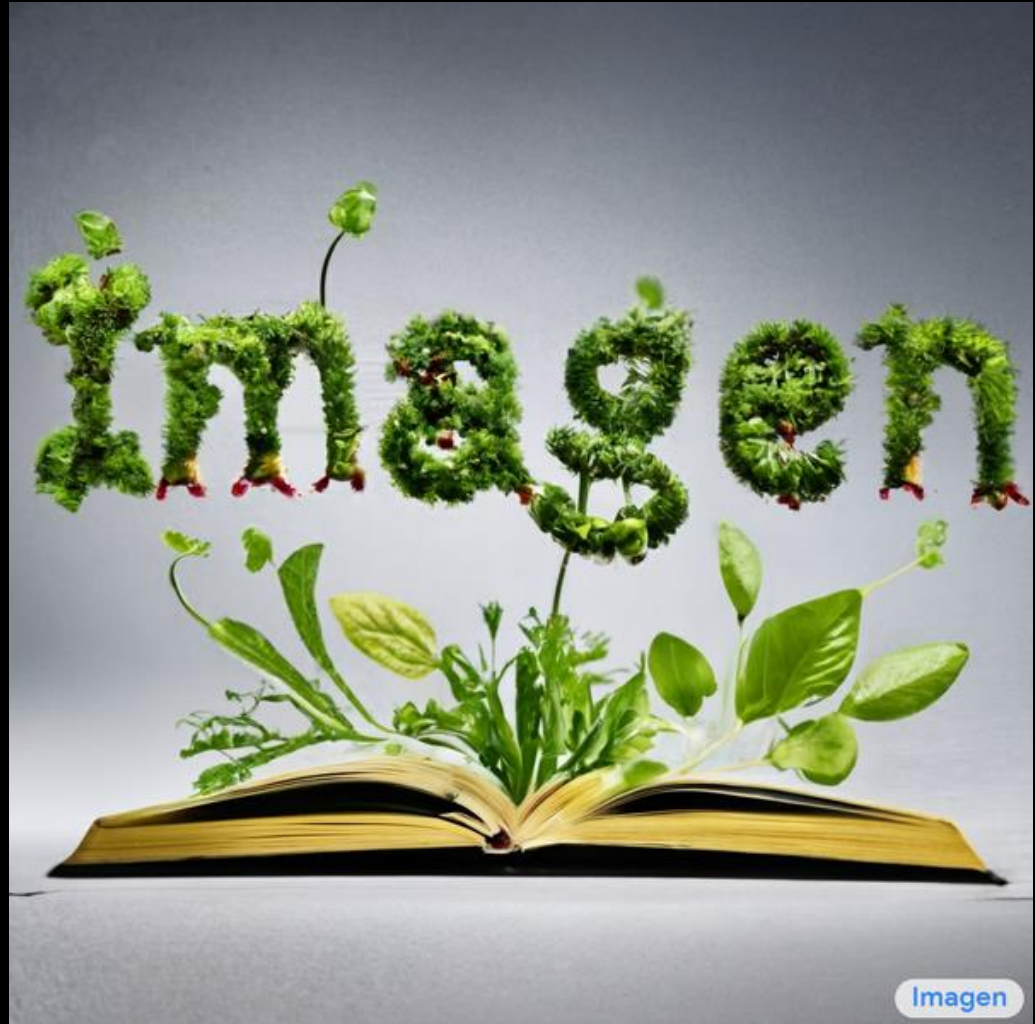**TranscribeGlass™**

Computer Assisted Real-time
Translation (CART)

# Existing

Imagen

Diffusion model + LLM(T5)

Text to image

# Existing

DALL-E 3

Diffusion model + LLM(GPT-4)


Text to image