

Vision and Text Integration

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
لعلوم والتكنولوجيا
King Abdullah University of
Science and Technology



LMH
Lady Margaret Hall

July 25, 2025

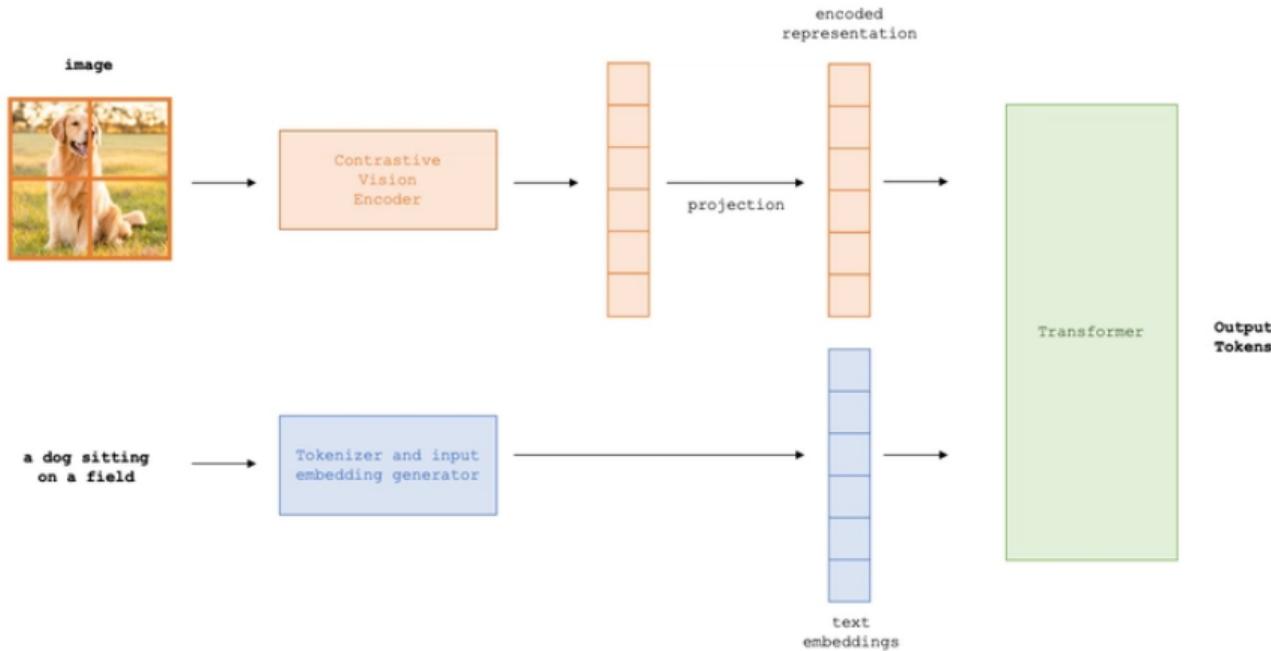


Table of Contents

1. Motivation
2. Learning Outcomes
3. Image and Video Captioning
 1. Image Captioning Pipeline
 2. Show and Tell
 3. Video Captioning
 4. Challenges in Captioning
 5. Recent Advances and Datasets
4. Multimodal Models: Image and Video Captioning
 1. Multimodal Representation Learning
 2. Joint Embedding Spaces for Image-Text Pairs
 3. Contrastive Learning for Multimodal Models

Table of Contents (cont.)

4. CLIP: Contrastive Language–Image Pretraining
5. Text-to-Image Generation Techniques
 1. Overview
 2. GAN-based Approaches
 3. Diffusion-based Approaches
6. Text-Conditioned Generation
 1. Overview of Text-Conditioned Generation
 2. Prompt Engineering
 3. Challenges in Conditioning
7. Classifier-Free Guidance (CFG)
8. Summary
9. References

- ▶ The human brain naturally integrates vision and language (e.g., describing scenes, reading captions).
- ▶ This inspires artificial systems to process and relate visual and textual information.
- ▶ Recent applications:
 - Image captioning
 - Visual question answering (VQA)
 - Generative AI

- ▶ Joint vision-language understanding is crucial for intelligent systems.
- ▶ Foundation for:
 - Autonomous agents
 - Assistive technologies
 - Content creation tools

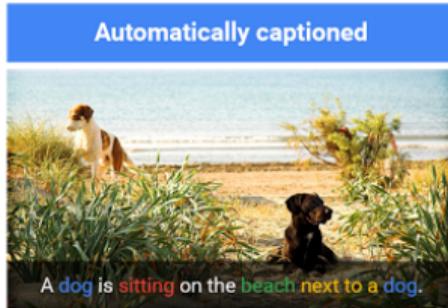
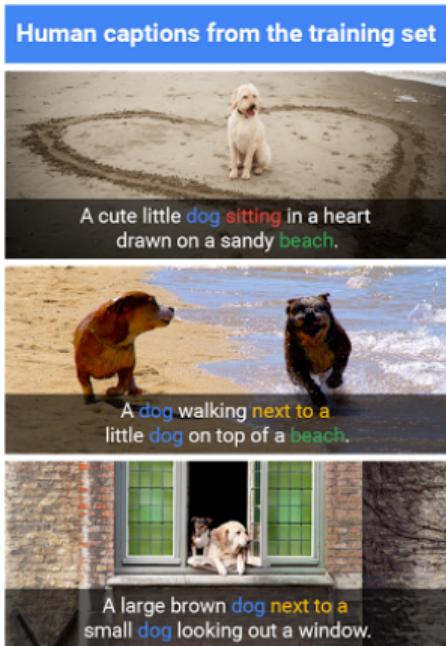
By the end of this session, learners will be able to:

- ▶ Understand how vision and text can be jointly modeled.
- ▶ Explore state-of-the-art models such as CLIP and diffusion-based text-to-image (T2I) methods.
- ▶ Implement and evaluate image and video captioning systems.
- ▶ Grasp classifier-free guidance and text-conditioned generation techniques.
- ▶ Analyze the limitations of current approaches and propose future research directions.

Vision and Text Integration: **Image and Video Captioning**

What is Captioning?

- ▶ **Definition:** Generating a descriptive text from an image or a video.



What is Captioning? (cont.)

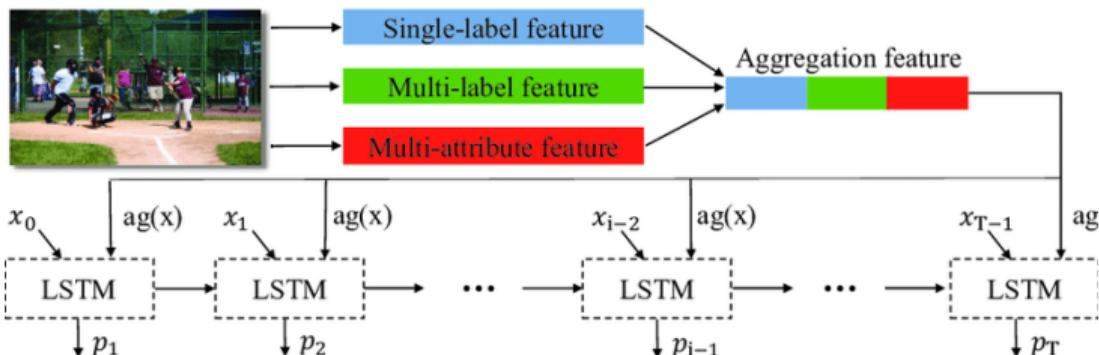
► Challenges:

- Semantics
- Scene understanding
- Temporal context

► Applications:

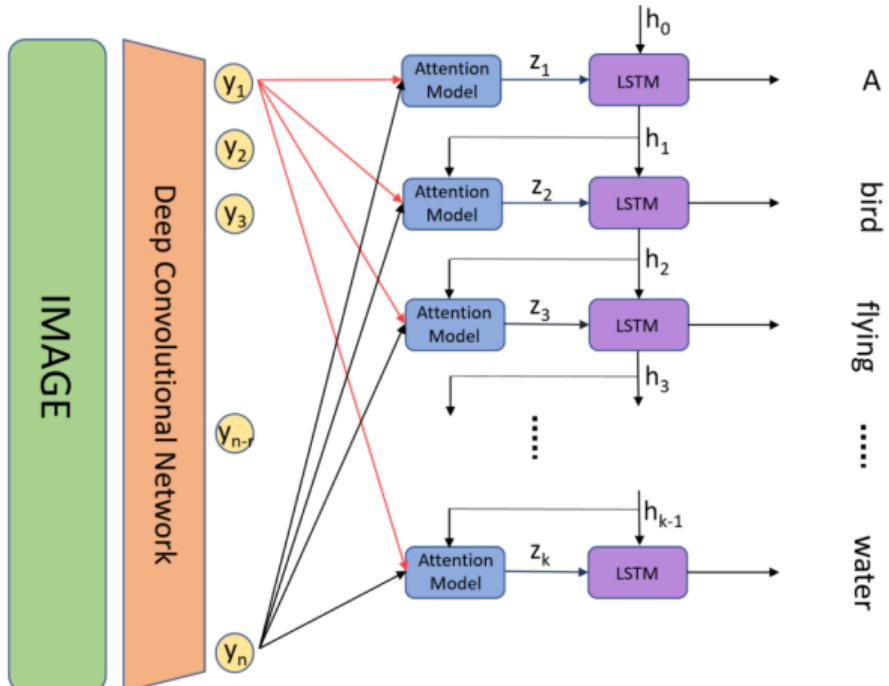
- Accessibility
- Content indexing
- Social media
- Surveillance

- ▶ **Feature Extraction:** Use a CNN (e.g., ResNet) to extract visual features from the image.
- ▶ **Language Modeling:** Use an RNN, LSTM, or Transformer to generate a textual description based on the extracted features.
- ▶ **Encoder-Decoder Architecture:** The CNN acts as the encoder and the RNN/LSTM/Transformer as the decoder.



Typical image captioning pipeline.

Image Captioning Pipeline (cont.)



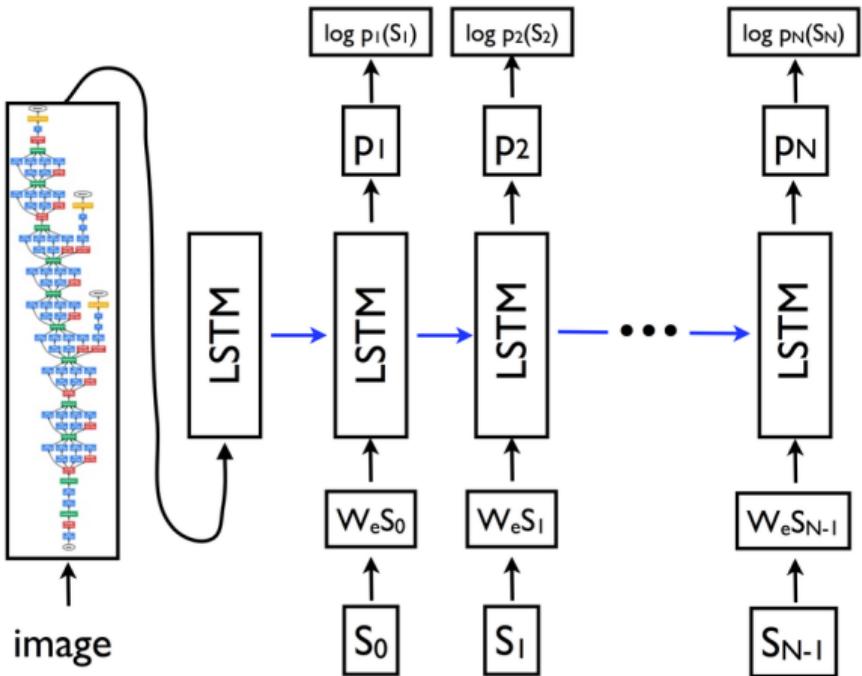
Active Learning in Image Captioning

Example: Show and Tell (Vinyals et al., 2015)



- ▶ One of the first successful end-to-end image captioning models.
- ▶ Uses a CNN (Inception) as encoder and an LSTM as decoder.
- ▶ Trained to maximize the likelihood of the correct caption given the image.

Example: Show and Tell (Vinyals et al., 2015) (cont)



Show and Tell model architecture (Vinyals et al., 2015).

Example: Show and Tell (Vinyals et al., 2015) (cont)

Helping visually impaired people better understand the content of images on the web.

- Pascal dataset
- Flickr30k dataset
- SBU dataset
- COCO dataset

$I \rightarrow$ input image

$S = \{S_1, S_2, \dots\} \rightarrow$ target sequence of words

$p(S|I) \rightarrow$ likelihood

S_t comes from a given dictionary

NIC: Neural Image Caption

$$\theta^* = \arg \max_{\theta} \sum_{(I, S)} \log p(S|I; \theta)$$

$$\log p(S|I) = \sum_{t=0}^N \log \underbrace{p(S_t|I, S_0, \dots, S_{t-1})}_{\text{model with an RNN}}$$

LSTM $h_{t+1} = f(h_t, x_t)$

CNN input (image & words)

$$i_t = \sigma(W_{ix}x_t + W_{im}m_{t-1})$$

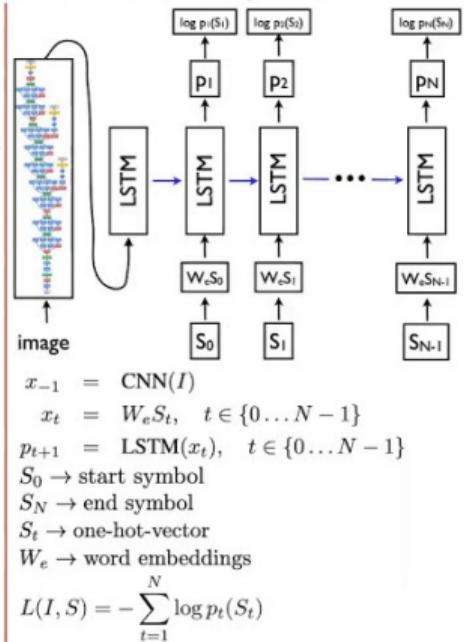
$$f_t = \sigma(W_{fx}x_t + W_{fm}m_{t-1})$$

$$o_t = \sigma(W_{ox}x_t + W_{om}m_{t-1})$$

$$c_t = f_t \odot c_{t-1} + i_t \odot h(W_{cx}x_t + W_{cm}m_{t-1})$$

$$m_t = o_t \odot c_t$$

$$p_{t+1} = \text{Softmax}(m_t)$$



BeamSearch: Iteratively consider the set of the k best sentences up to time t as candidates to generate sentences of size $t + 1$, and keep only the resulting best k of them.
(beam size: $k = 20$)



Describes without errors Describes with minor errors Somewhat related to the image Completely related to the image

Vinyals, Oriol, et al. "Show and tell: A neural image caption generator." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.

Example: Show and Tell (Vinyals et al., 2015) (cont)



Describes without errors

Describes with minor errors

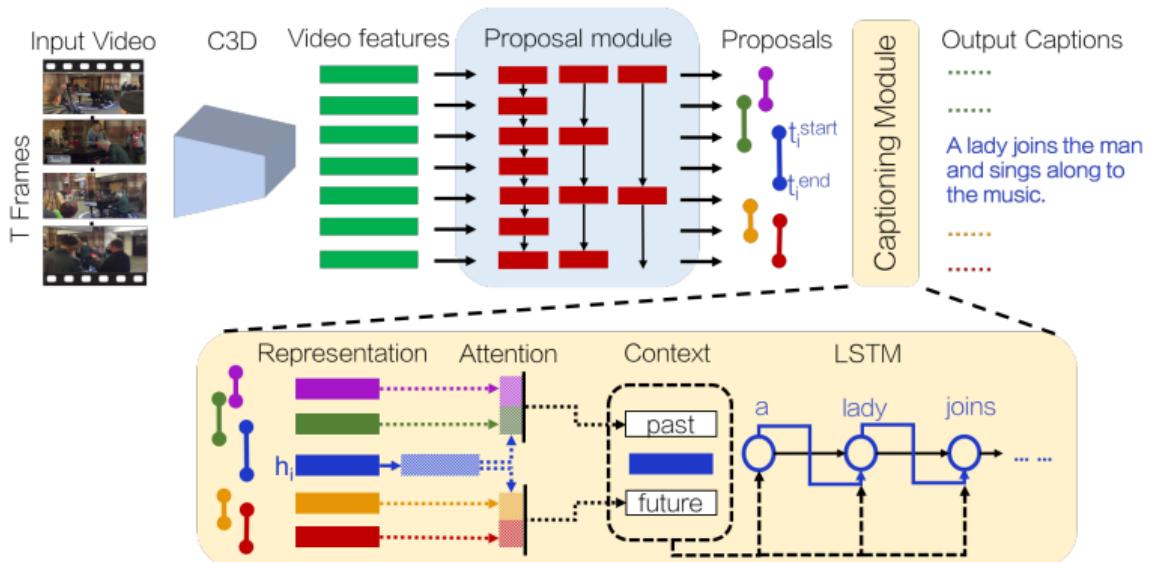
Somewhat related to the image

Unrelated to the image

Figure 5. A selection of evaluation results, grouped by human rating.

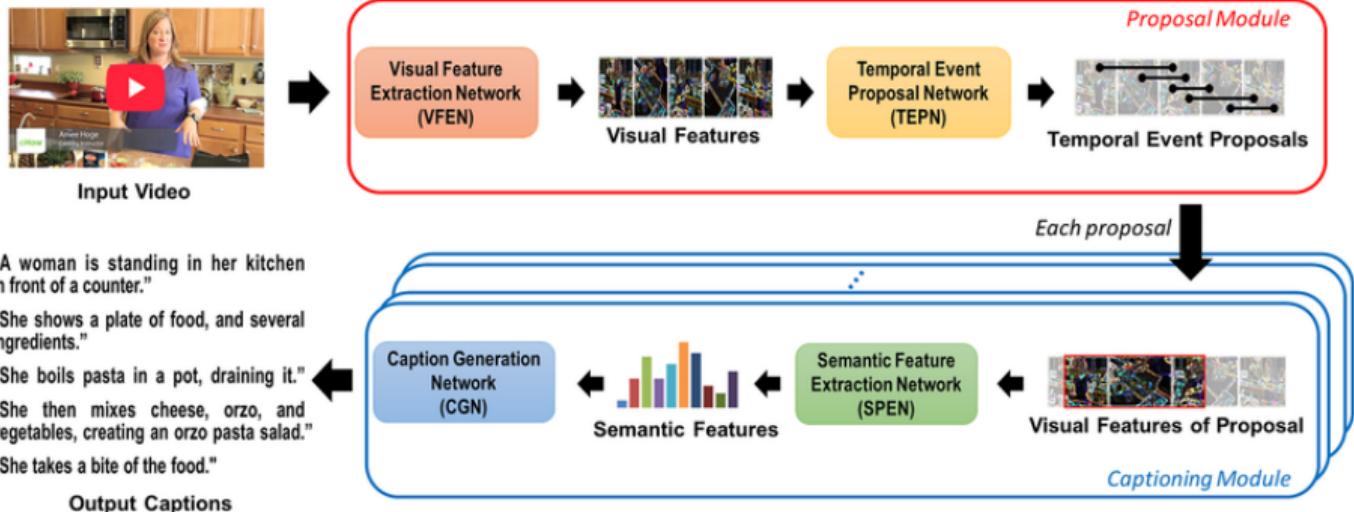
- ▶ **Video Captioning:** Generating descriptive text for a sequence of video frames.
- ▶ **Temporal Sequence Modeling:** Adds modeling of temporal dependencies between frames, in addition to spatial features.
- ▶ **Common Architecture:**
 - CNN for frame-level feature extraction.
 - RNN (e.g., LSTM/GRU) for modeling temporal sequence.
 - Attention mechanisms to focus on relevant frames or regions.
- ▶ **Recent Advances:**
 - Transformer-based models for parallel sequence modeling.
 - Memory networks for capturing long-term dependencies.

Video Captioning (cont.)



Dense-Captioning Events in Videos (Hendricks et al., 2017)

Video Captioning (cont.)



Video Captioning (cont.)



"A woman is standing in her kitchen in front of a counter."



"She shows a plate of food, and several ingredients."



"She boils pasta in a pot, draining it."



"She then mixes cheese, orzo, and vegetables, creating an orzo pasta salad."



"She takes a bite of the food."

- ▶ **Ambiguity:** Multiple valid captions for the same image or video.
- ▶ **Context Understanding:** Captions need to capture context and semantics accurately.
- ▶ **Temporal Dynamics:** In video captioning, capturing the dynamics of actions and events over time is crucial.
- ▶ **Data Scarcity:** High-quality annotated datasets for training captioning models are limited.
- ▶ **Evaluation Metrics:** Evaluating caption quality is subjective and often relies on metrics like BLEU, METEOR, or CIDEr, which may not align with human judgment.

Recent Advances in Captioning

Image Captioning Advances: OSCAR: Object-Semantics Aligned Pre-training

for Vision-and-Language Tasks.

OSCaR: Object State Captioning and State Change Representation

Nguyen Nguyen¹, Jing Bi¹, Ali Vosoughi¹, Yapeng Tian², Pooyan Fazli³, Chenliang Xu¹

¹University of Rochester, ²University of Texas at Dallas, ³Arizona State University

{nguyen.nguyen, jing.bi, ali.vosoughi, chenliang.xu}@rochester.edu,
yapeng.tian@utdallas.edu, pooyan@asu.edu

Abstract

The capability of intelligent models to extrapolate and comprehend changes in object states is a crucial yet demanding aspect of AI research, particularly through the lens of human interaction in real-world settings. This task involves describing complex visual environments, identifying active objects, and interpreting their changes as conveyed through language. Traditional methods, which isolate object captioning and state change detection, offer a limited view of dynamic environments. Moreover, relying on a small set of symbolic words to represent changes has restricted the expressiveness of language. To address these challenges, in this paper, we introduce the Object State Captioning and State Change Representation (OSCaR) dataset and benchmark. OSCaR consists of 14,084 annotated video segments with nearly

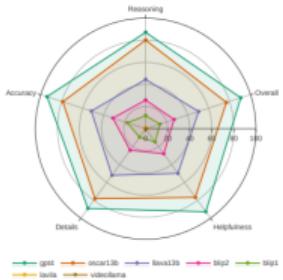
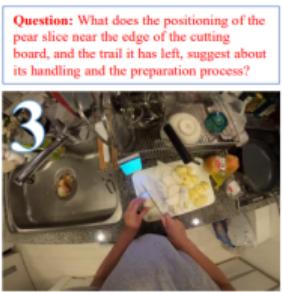
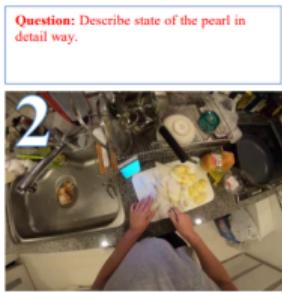
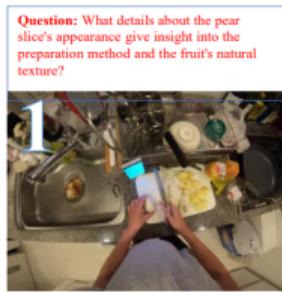
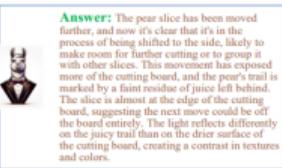
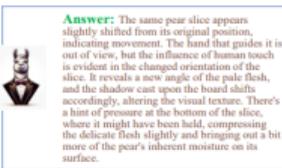
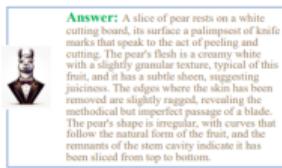


Figure 1: Surpassing prior models in aligning with human judgements. Our method achieves near parity with GPT-4V ratings across helpfulness, accuracy, reasoning, and other key metrics.

OSCAR: Object-Semantics Aligned Pre-training for Vision-and-Language Tasks.



Question: What does the sequence of movements of the pear slice across the cutting board indicate about the method of preparation?



Answer: Over the sequence of these three frames, the pear slice undergoes a subtle but telling transformation through human intervention. Initially at rest, it is then picked up and moved, each action causing a change in its position and the way it interacts with its environment. This movement is both purposeful, and the slight compressions and shifts reveal not just the hand's manipulation, but also the pear's textural responses to being handled. The progression shows a choreography of culinary preparation, where each movement is part of a larger, rhythmic and method to the task at hand. The final frame indicates the continuation of this process beyond the captured moment, as the slice approaches the edge of the board, ready for the next step in its journey from whole fruit to ingredient.

Figure 2: OSCaR's description of state, state change, and illustration of reasoning. State description involves the characterization of a specific region of interest within the video and the associated activity. State change entails the description of the evolution of a system over a defined temporal sequence. Furthermore, the analysis of the state of an object is centered on comprehending and elucidating the mechanisms underlying the object's evolution.

VinVL: Large-scale pre-training with improved object-visual representations.

VinVL: Revisiting Visual Representations in Vision-Language Models

Pengchuan Zhang^{♡†}

Xiujun Li^{♡♣†}

Xiaowei Hu[♡]

Jianwei Yang[♡]

Lei Zhang[♡]

Lijuan Wang[♡]

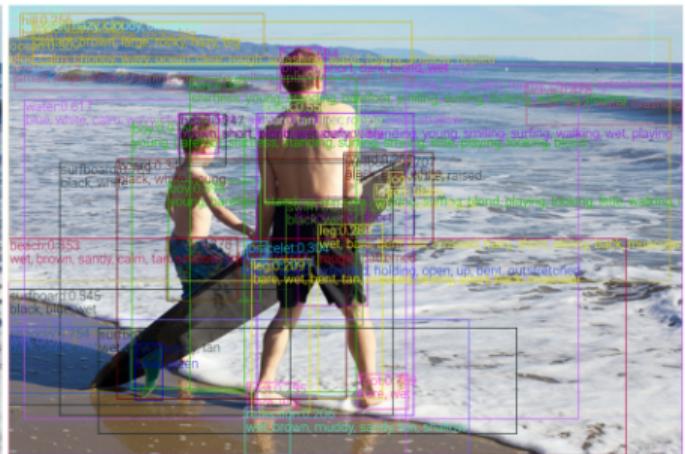
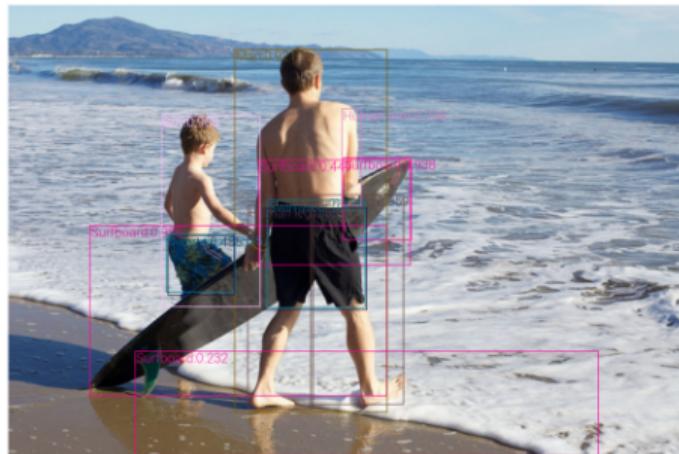
Yejin Choi[♣]

Jianfeng Gao[♡]

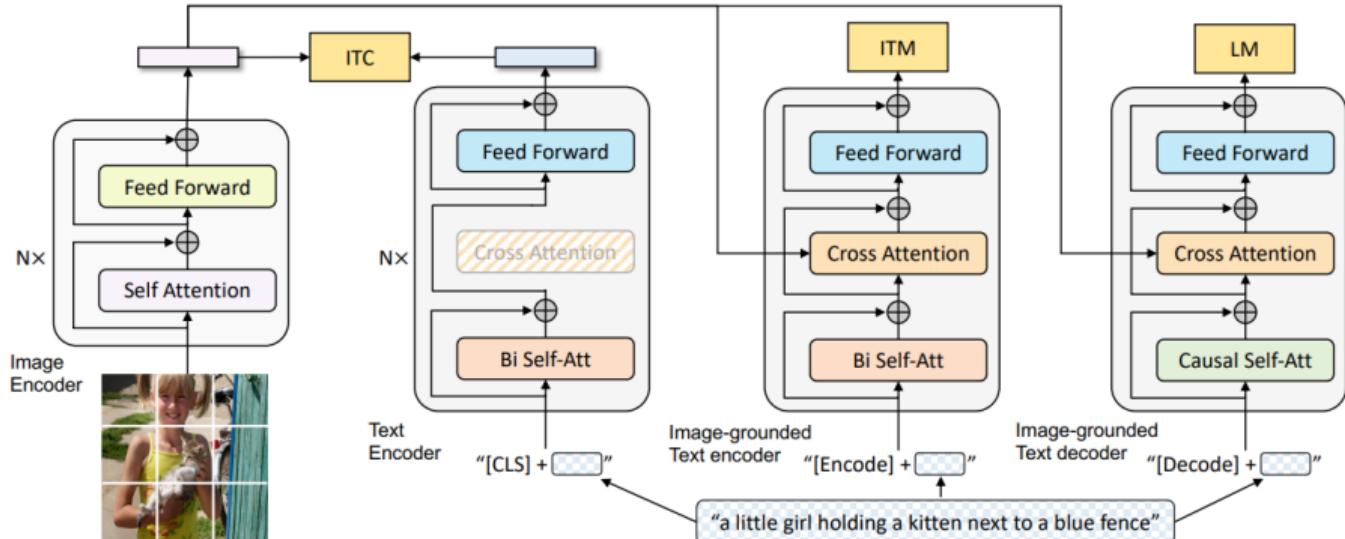
March 11, 2021

Recent Advances in Captioning (cont.)

VinVL: Large-scale pre-training with improved object-visual representations.



BLIP: Bootstrapped Language-Image Pretraining for unified vision-language understanding and generation.



Video Captioning Advances:

- ▶ **MARN:** Multi-modal Attentive Recurrent Network for video captioning.
- ▶ **STG-KD:** Spatio-Temporal Graph Knowledge Distillation for video captioning.
- ▶ **ClipCap:** Zero-shot image-to-text generation using CLIP and language models.

Key Datasets

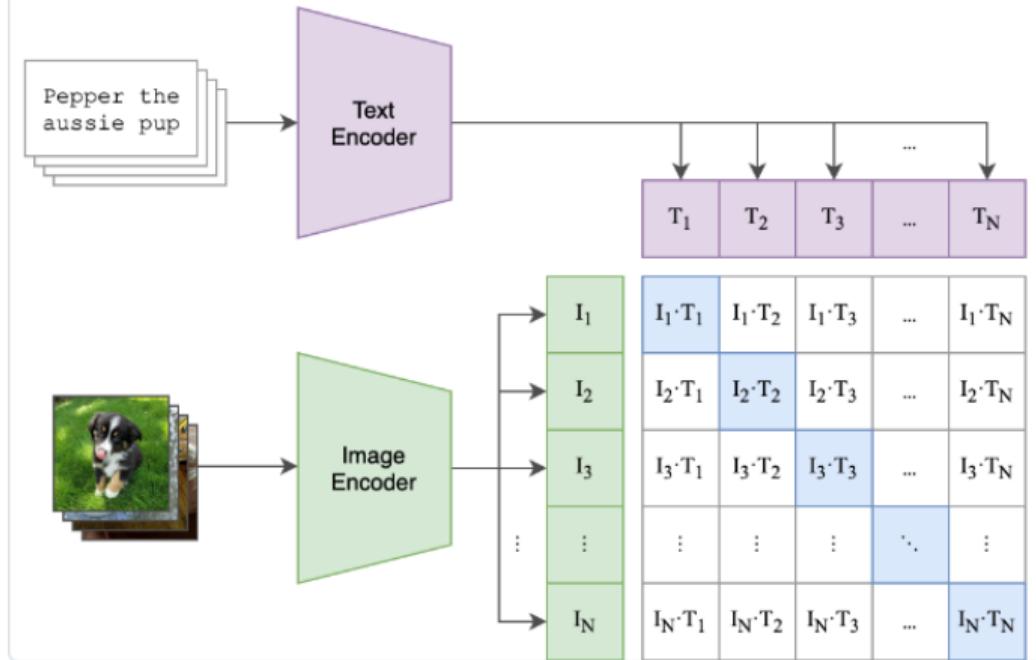
- ▶ **MSCOCO:** Large-scale dataset for image captioning.
- ▶ **ActivityNet Captions:** Benchmark for dense video captioning.
- ▶ **YouCook2:** Large-scale dataset of instructional cooking videos with captions.

Image and Video Captioning: **Multimodal Models**

- ▶ **Definition:** Multimodal models integrate visual and textual information to generate descriptive captions for images and videos.
- ▶ **Challenges:**
 - Ambiguity in visual content
 - Context understanding
 - Temporal dynamics in videos
- ▶ **Applications:**
 - Image captioning
 - Video captioning
 - Assistive technologies (e.g., for visually impaired)

- ▶ **Goal:** Learn shared representations for different modalities (e.g., images and text).
- ▶ **Approach:** Map both images and text into a joint embedding space.
- ▶ **Benefit:** Enables comparison and retrieval across modalities.

Multimodal Representation Learning (cont.)



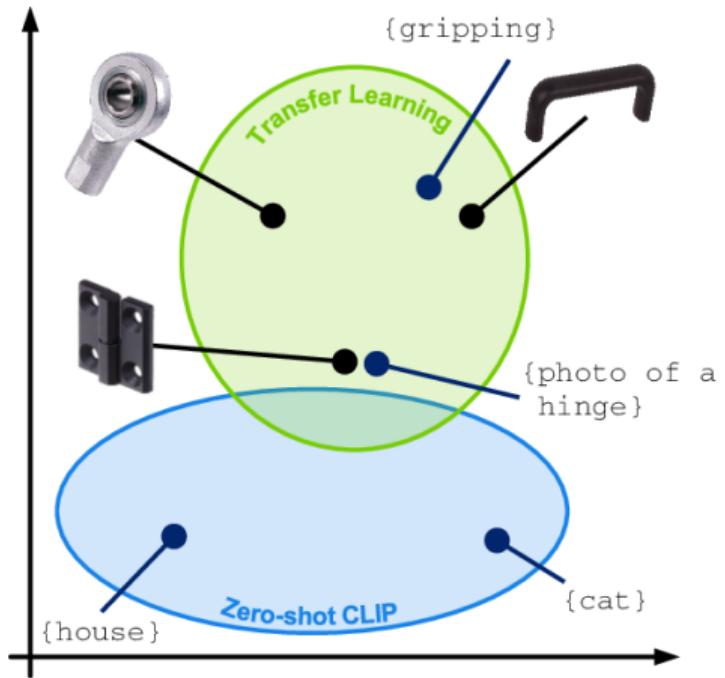
Multimodal representation learning: mapping images and text into a joint embedding space.

Joint Embedding Spaces for Image-Text Pairs



- ▶ Images and their corresponding captions are projected into a common vector space.
- ▶ Matching image-text pairs are close together; non-matching pairs are far apart.
- ▶ Used in tasks like cross-modal retrieval and caption generation.

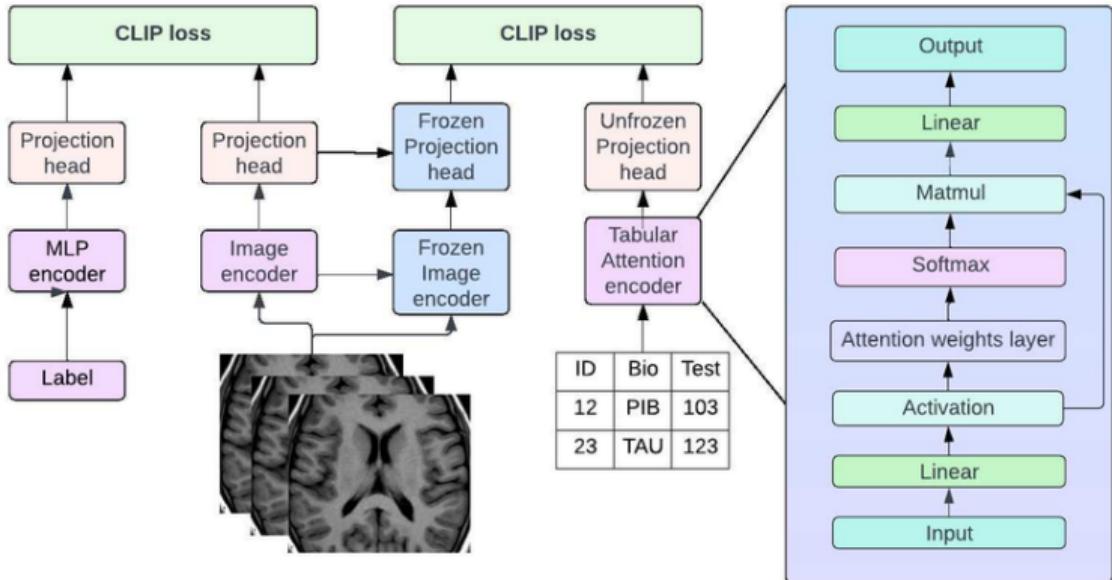
Joint Embedding Spaces for Image-Text Pairs (cont.)



Joint embedding space for image-text pairs: matching pairs are close, non-matching pairs are distant.

- ▶ **Contrastive Objective:** Pull matching image-text pairs together, push non-matching pairs apart.
- ▶ **Popular Methods:** CLIP, ALIGN, etc.
- ▶ **Loss Function:** Often uses InfoNCE or similar contrastive losses.

Contrastive Learning for Multimodal Models (cont.)

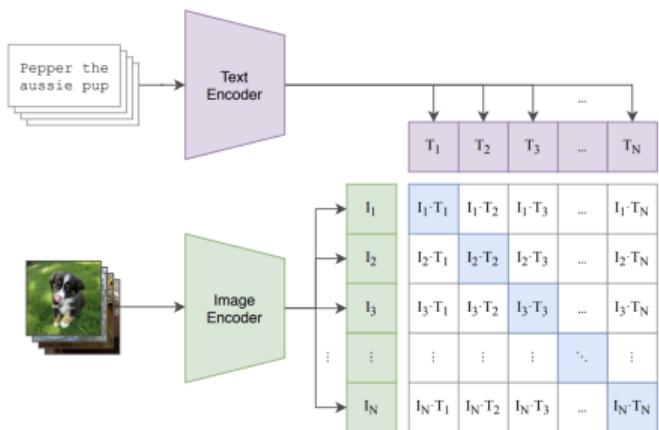


Contrastive learning for multimodal models: aligning image-text pairs in a joint embedding space.

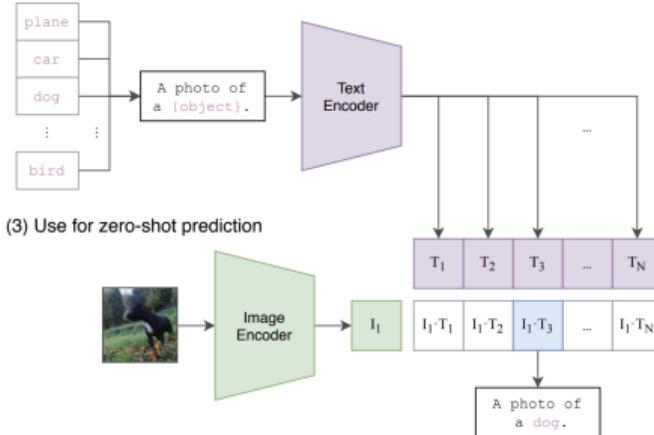
- ▶ **Overview:** CLIP is a multimodal model trained to connect images and text using contrastive learning.
- ▶ **Training Data:** 400 million image-text pairs collected from the internet.
- ▶ **Architecture:**
 - Separate Transformer encoders for images and text.
 - Both modalities are mapped into a shared embedding space.
- ▶ **Contrastive Loss:** Uses dot-product similarity to align matching image-text pairs and separate non-matching pairs.
- ▶ **Zero-shot Classification:** Enables classification of images without task-specific training by comparing image embeddings to text prompt embeddings.

CLIP (Contrastive Language–Image Pretraining) (cont.)

(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

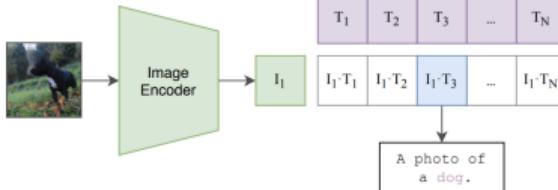


Figure 1. Summary of our approach. While standard image models jointly train an image feature extractor and a linear classifier to predict some label, CLIP jointly trains an image encoder and a text encoder to predict the correct pairings of a batch of (image, text) training examples. At test time the learned text encoder synthesizes a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes.

CLIP architecture: separate encoders for images and text, mapping both into a shared embedding space.

Why CLIP Matters:

- ▶ **Open-vocabulary vision models:** CLIP can recognize and describe a wide range of concepts beyond fixed label sets.
- ▶ **Robust to distribution shift:** Performs well on images from domains not seen during training.
- ▶ **No task-specific training required:** Enables zero-shot learning for new tasks by leveraging text prompts.

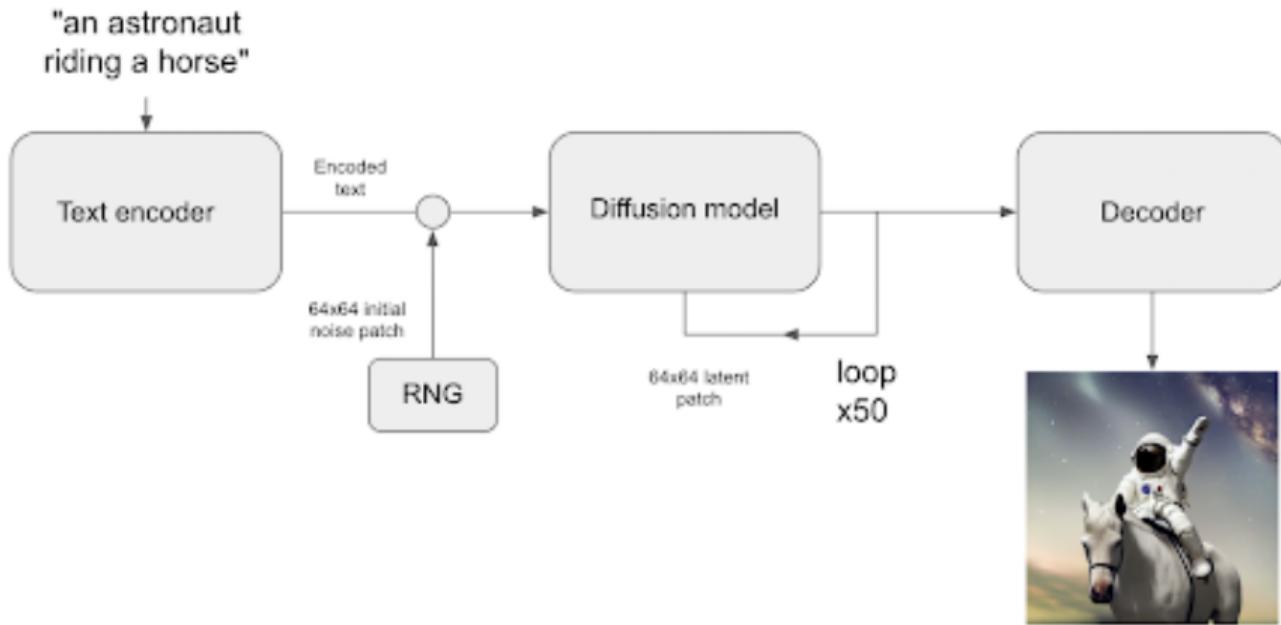
Limitations of CLIP:

- ▶ **Struggles with fine-grained details:** May fail to distinguish subtle differences between visually similar objects.
- ▶ **Lacks true reasoning abilities:** Cannot perform complex reasoning or understand relationships beyond surface-level associations.
- ▶ **Prone to dataset biases:** Inherits biases present in the large-scale web data used for training.

Text-to-Image Generation Techniques

- ▶ Text-to-image generation involves creating images from textual descriptions.
- ▶ Key techniques include:
 - Diffusion models
 - Generative adversarial networks (GANs)
 - Variational autoencoders (VAEs)
- ▶ Applications:
 - Art generation
 - Content creation
 - Virtual environments

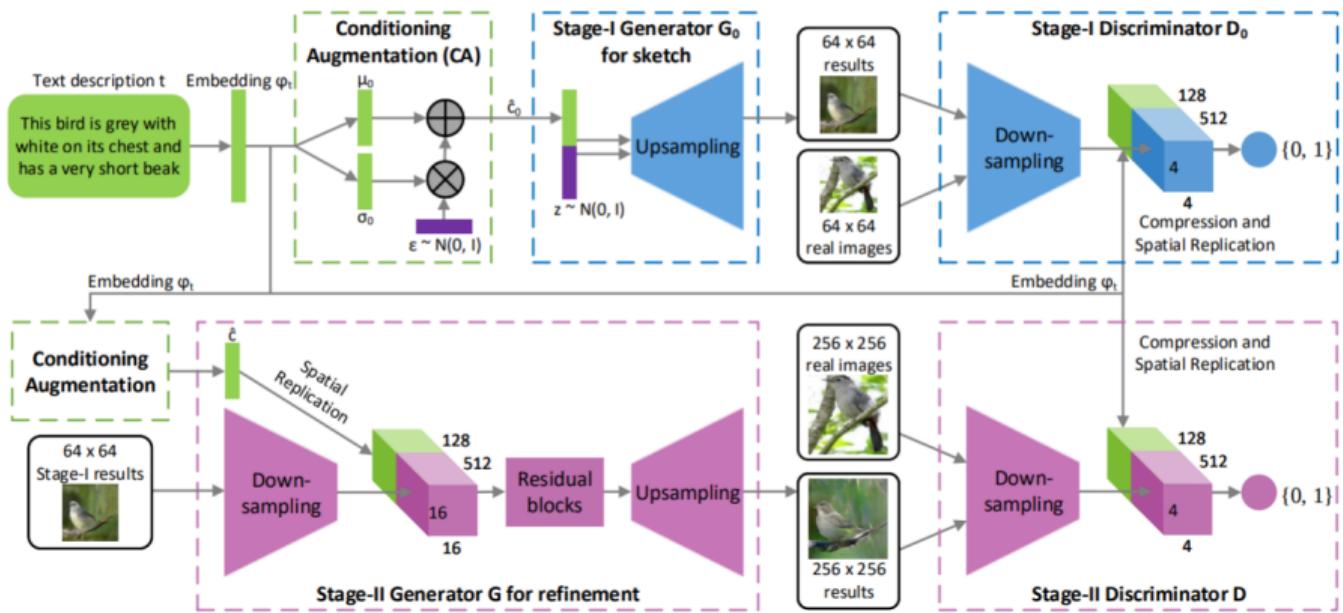
Text-to-Image Generation (cont.)



- ▶ **Goal:** Generate realistic images conditioned on text prompts.
- ▶ Requires alignment between vision and language features.
- ▶ **Applications:**
 - Design
 - Entertainment
 - Accessibility

GAN-based Approaches

StackGAN: Generates images in multiple stages, from coarse to fine details.



GAN-based Approaches (cont.)

Text description	This bird is blue with white and has a very short beak	This bird has wings that are brown and has a yellow belly	A white bird with a black crown and yellow beak	This bird is white, black, and brown in color, with a brown beak	The bird has small beak, with reddish brown crown and gray belly	This is a small, black bird with a white breast and white on the wingbars.	This bird is white black and yellow in color, with a short black beak
Stage-I images							
Stage-II images							

GAN-based Approaches (cont.)

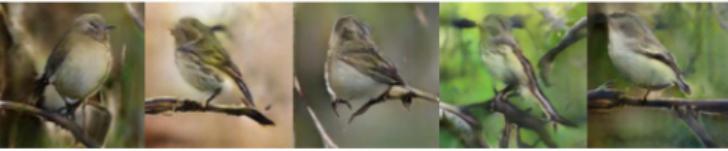
This small blue bird has a short pointy beak and brown on its wings



This bird is completely red with black wings and pointy beak



A small sized bird that has a cream belly and a short pointed bill



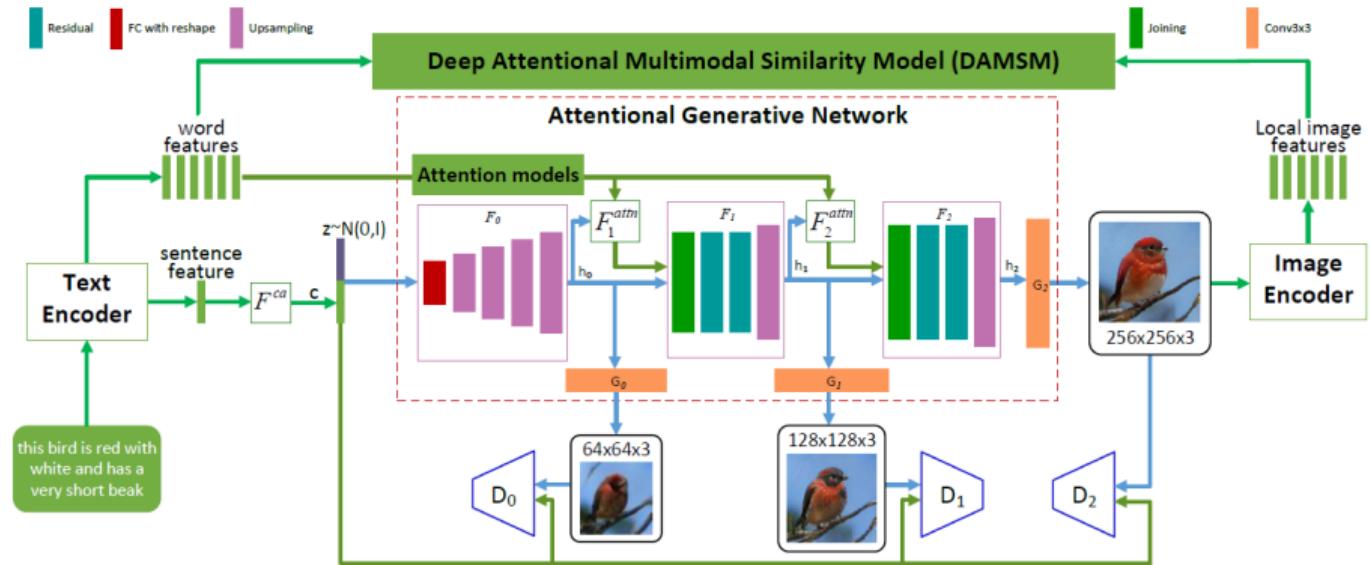
A small bird with a black head and wings and features grey wings



Figure 7. Birds with different poses and viewpoints generated with the same input text embedding by our StackGAN. The noise vector z and text embedding are fixed for each row.

GAN-based Approaches (cont.)

AttnGAN: Incorporates attention mechanisms to better align image regions with words in the text.



GAN-based Approaches (cont.)

this bird is red with white and has a very short beak



Figure 1. Example results of the proposed AttnGAN. The first row gives the low-to-high resolution images generated by G_0 , G_1 and G_2 of the AttnGAN; the second and third row shows the top-5 most attended words by F_1^{attn} and F_2^{attn} of the AttnGAN, respectively. Here, images of G_0 and G_1 are bilinearly upsampled to have the same size as that of G_2 for better visualization.

Stage-wise Generation:

- ▶ Images are generated progressively, refining details at each stage.

Limitations:

- ▶ Training instability
- ▶ Mode collapse (lack of diversity in generated images)

- ▶ Diffusion models generate images by iteratively denoising random noise, conditioned on text prompts.
- ▶ Notable models:
 - DALL•E 2
 - Imagen
 - Stable Diffusion
- ▶ **Key idea:** Learn to reverse a gradual noising process, producing high-quality images from pure noise.
- ▶ **Advantages:**
 - Superior image quality and diversity compared to GANs and VAEs
 - Better alignment with textual descriptions

DALL•E 2: A diffusion model that generates images from text prompts, achieving high fidelity and diversity.

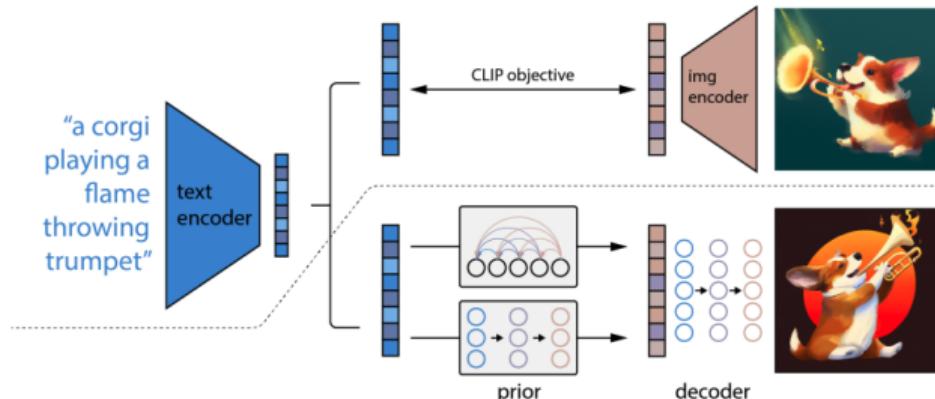


Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process, through which we learn a joint representation space for text and images. Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive or diffusion prior to produce an image embedding, and then this embedding is used to condition a diffusion decoder which produces a final image. Note that the CLIP model is frozen during training of the prior and decoder.

DALL•E 2 architecture (source: OpenAI)

Diffusion-based Approaches (cont.)

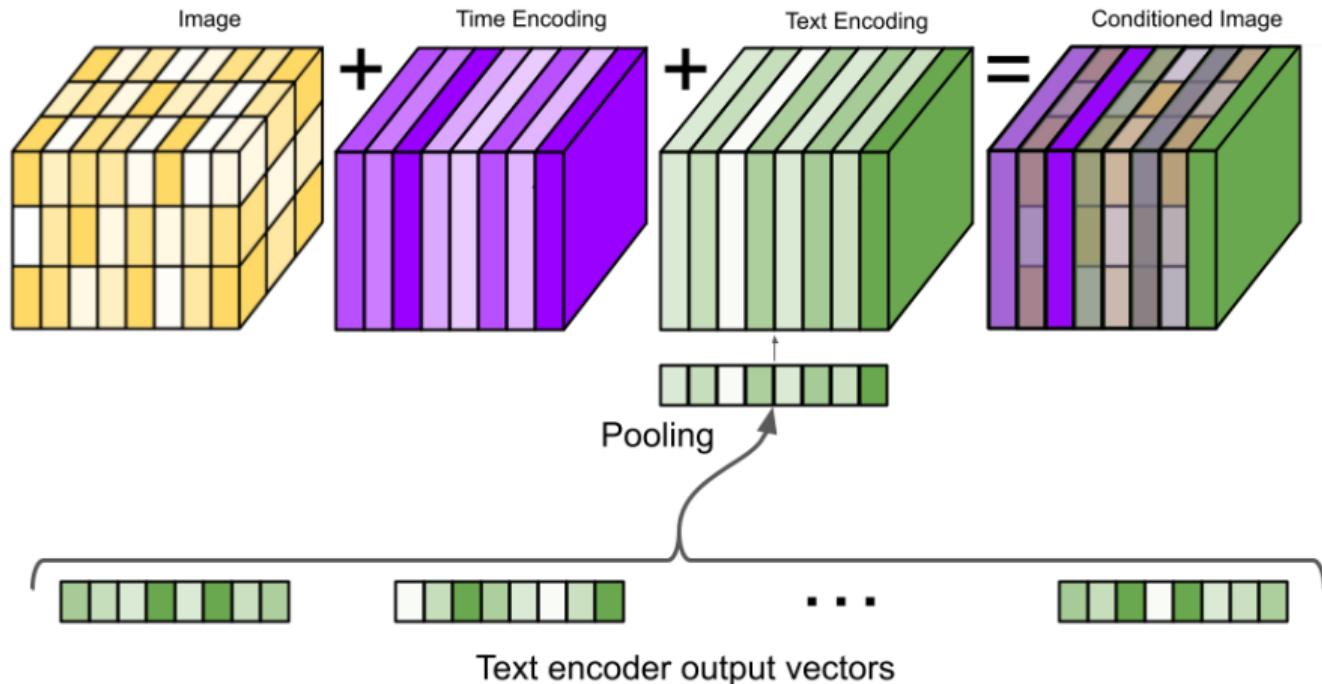


Diffusion-based Approaches (cont.)



Imagen: A diffusion model by Google that generates high-quality images from text, focusing on photorealism and fine details.

Diffusion-based Approaches (cont.)



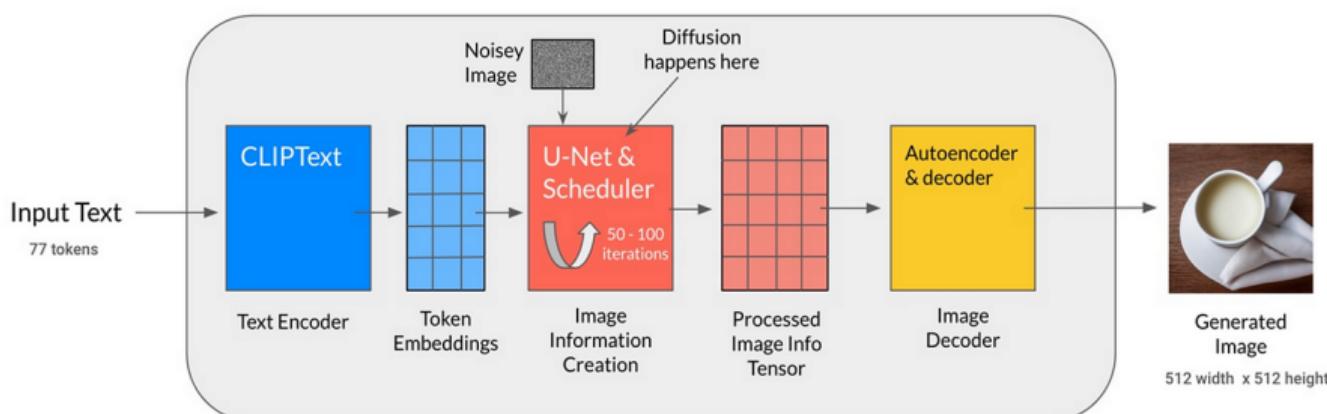
How Imagen works: from text to image generation (source: Google)

Diffusion-based Approaches (cont.)



Stable Diffusion: An open-source diffusion model that generates high-quality images from text, designed for efficiency and accessibility.

Stable Diffusion Architecture



Stable Diffusion architecture (source: Stability AI)

Diffusion-based Approaches (cont.)



Diffusion Model Architecture:

- ▶ **U-Net backbone:** The core neural network used for denoising, featuring an encoder-decoder structure with skip connections.
- ▶ **Text encoder:** Converts text prompts into embeddings. Common choices include T5 or CLIP Text Encoder.
- ▶ **Cross-attention:** Mechanism to inject text conditioning into the image generation process, allowing the model to align image features with textual information.

Text-Conditioned Generation

- ▶ Text-to-image generation involves creating images from textual descriptions.
- ▶ Applications include art generation, content creation, and visual storytelling.
- ▶ Key challenges:
 - Ensuring high fidelity to the text description.
 - Maintaining diversity in generated images.
 - Handling complex and abstract descriptions.

▶ Concatenation of Embeddings:

- Combine text and image embeddings before feeding them into the generative model.
- Simple and effective for early fusion of modalities.

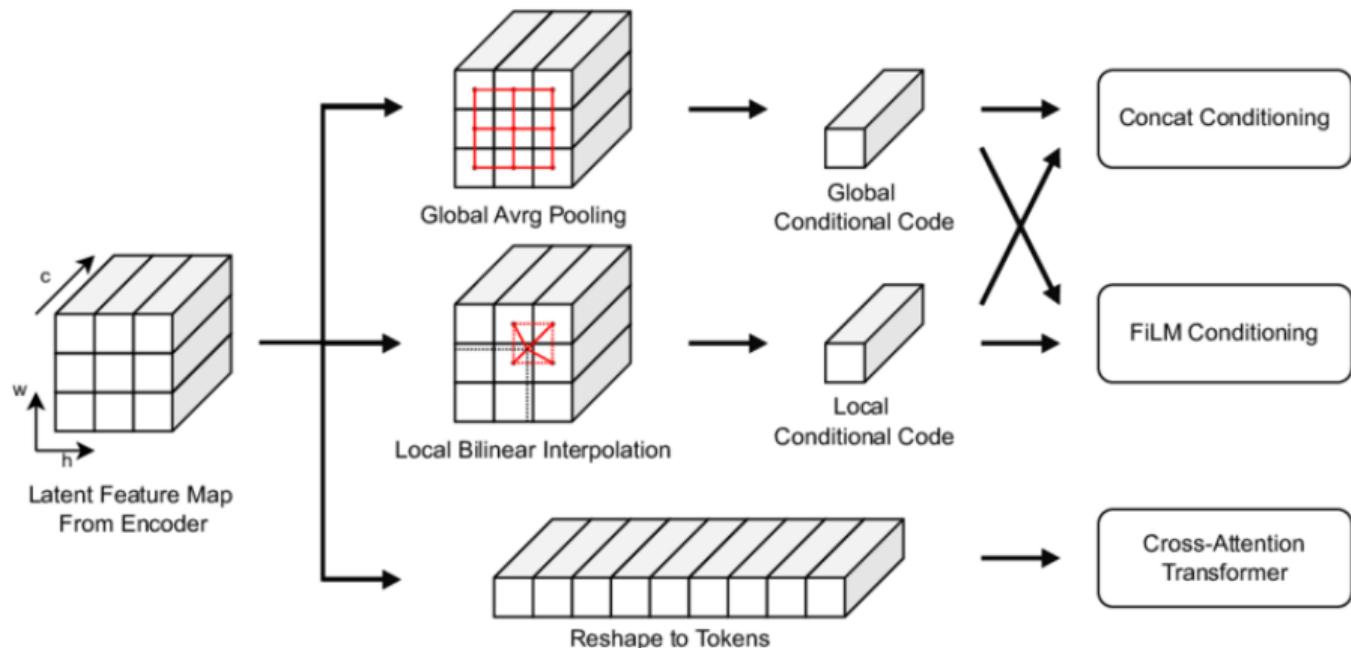
▶ Cross-Attention in Transformer Layers:

- Use cross-attention mechanisms to allow the model to focus on relevant parts of the text while generating images.
- Enables fine-grained alignment between text and image features.

▶ Classifier-Free Conditioning:

- Train the model with and without conditioning information.
- Allows flexible control over the influence of text during generation.

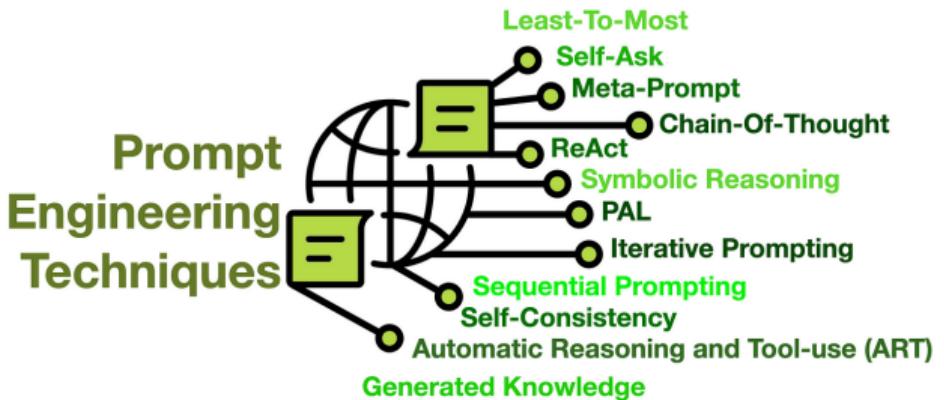
Text Conditioning Strategies (cont.)



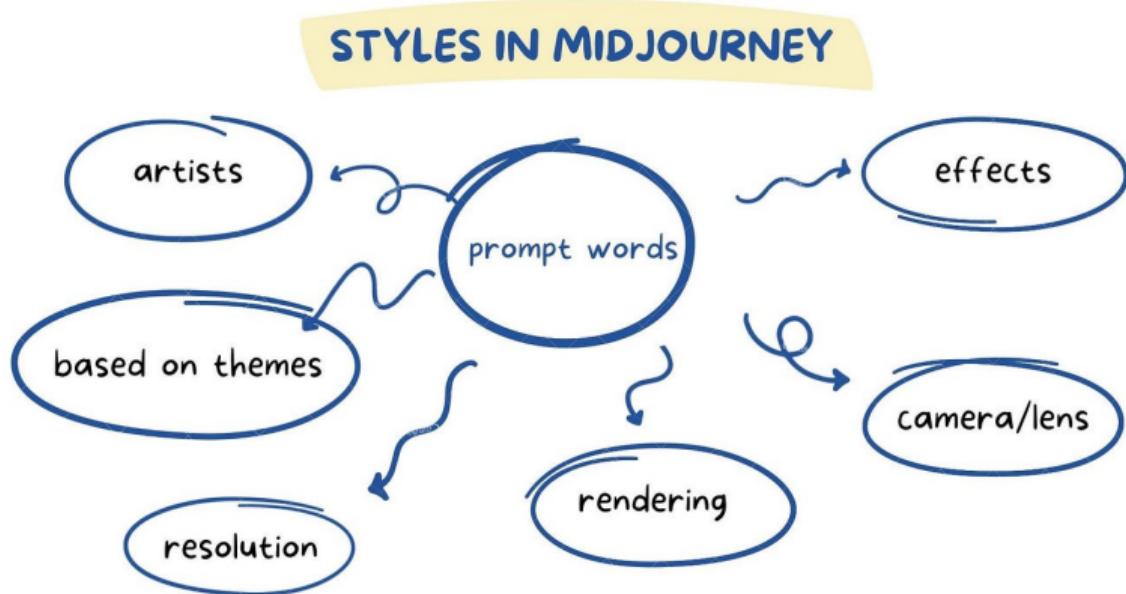
A visualization of text conditioning strategies in text-to-image generation.

- ▶ Crafting effective prompts is crucial for controlling the output of text-to-image models.
- ▶ Well-designed prompts can improve image quality, relevance, and adherence to user intent.
- ▶ Key aspects:
 - **Style Control:** Specify artistic styles (e.g., “in the style of Van Gogh”).
 - **Object Presence:** Clearly mention desired objects and their attributes.
 - **Layout:** Indicate spatial relationships (e.g., “a cat sitting on a chair”).

12 Prompt Engineering Techniques



- ▶ **Midjourney:** “A futuristic cityscape at sunset, vibrant colors, ultra-detailed, digital art”
- ▶ **DALL-E:** “An armchair in the shape of an avocado”
- ▶ **SDXL:** “A portrait of a medieval knight, oil painting, dramatic lighting”
- ▶ Experimenting with prompt phrasing can lead to diverse and creative outputs.



An advanced guide to writing prompts for Midjourney.

- ▶ **Ambiguity in Natural Language:** Text descriptions can be vague or open to multiple interpretations, making it difficult for models to generate the intended image.
- ▶ **Alignment with Low-Level Image Details:** Ensuring that fine-grained details in the image match the textual description remains a significant challenge.
- ▶ **Long Prompts and Hallucinations:** Handling lengthy or complex prompts can lead to hallucinated content or loss of important details in the generated images.



- | | |
|---|---|
| Ambiguous Prompts Lead to Unfocused Responses | • Be Specific |
| Clichés | • Example-based Prompts |
| Complexity Handling | • Chain of Thought Prompting |
| Inconsistent Tone or Style | • Persona-driven Prompts |
| Loading the Model with a lot of Context | • Maintain Balance and Specificity of Prompts |
| Risk i.e. Generating Untrue Information | • Augmented Generation |
| Data Privacy Issue | • Implement Robust Data Governance |
| Iterative Refining of Prompts | • Embrace Iterative Refinement |

www.gsdccouncil.org

Some Prompt Engineering Challenges and Their Solutions.

Classifier-Free Guidance (CFG)

What is Classifier-Free Guidance?

- ▶ A technique to improve the quality of text-to-image generation by controlling the influence of text prompts.
- ▶ It allows for generating images that are more aligned with the provided text while maintaining diversity.
- ▶ CFG is particularly useful in diffusion models, enhancing the coherence between text and generated images.

- ▶ **Balances realism and fidelity:** Adjusts the trade-off between generating realistic images and staying faithful to the text prompt.
- ▶ **No external classifier needed:** Removes the requirement for a separate classifier to guide the generation process.
- ▶ **Score interpolation:** Interpolates between conditional (with prompt) and unconditional (without prompt) model outputs to control guidance strength.

Key Properties of Classifier-Free Guidance (cont.)

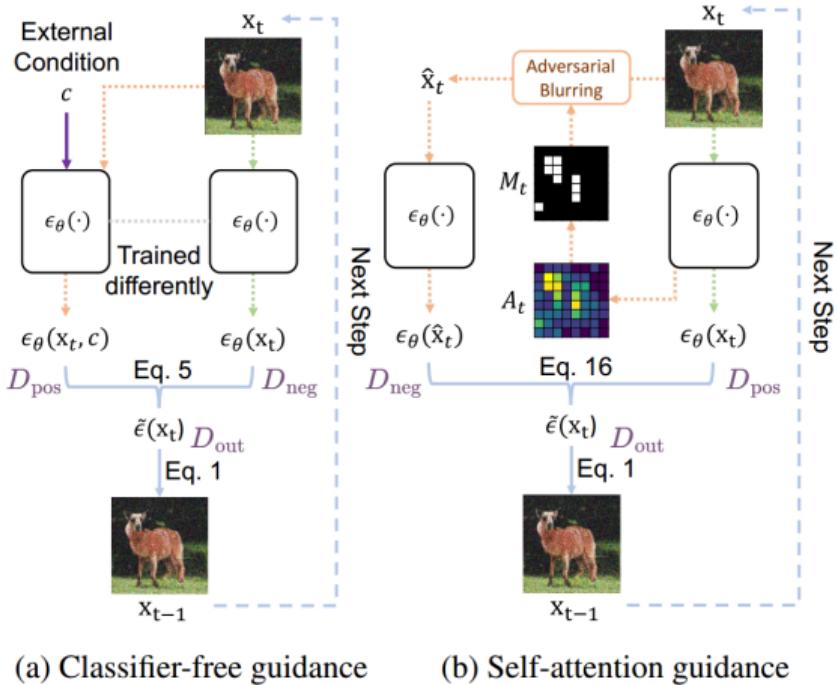


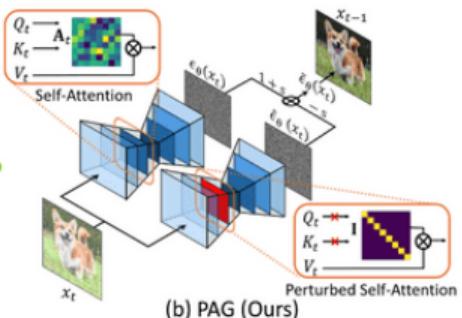
Illustration of Classifier-Free Guidance in text-to-image generation.

How does it work?

- ▶ The model generates images based on both the text prompt and random noise.
- ▶ By adjusting the guidance scale, users can control how much the text influences the final image.
- ▶ Higher guidance scales lead to images that are more closely aligned with the text, while lower scales allow for more creative freedom.

An overview of CLASSIFIER-FREE DIFFUSION GUIDANCE methods

- Training-free?
- Applicable to unconditional models?
- Single model?
- Train separate impaired model?



$$\hat{D}_{\text{out}}(x | \sigma) = D_{\text{neg}}(x | \sigma) + (1 + \gamma) (D_{\text{pos}}(x | \sigma) - D_{\text{neg}}(x | \sigma)).$$

Illustration of Classifier-Free Guidance in text-to-image generation.

- ▶ **Duplicate forward pass:** For each input, perform two forward passes through the model:

- One with the conditioning (e.g., text prompt)
- One without conditioning (unconditional)

- ▶ **Combine outputs:** The final score is computed as:

$$s = s_{\text{uncond}} + w \cdot (s_{\text{cond}} - s_{\text{uncond}})$$

where:

- s_{cond} is the model output with conditioning
- s_{uncond} is the model output without conditioning
- w is the guidance scale parameter

- ▶ **Effect of w :** Controls the strength of guidance. Higher w enforces stronger adherence to the prompt.

- ▶ **Simpler and faster:** Eliminates the need for a separate classifier, reducing complexity and computational overhead.
- ▶ **Effective for text-guided diffusion:** Provides strong alignment between generated images and text prompts in diffusion models.
- ▶ **Foundation of modern T2I systems:** Forms the basis of Stable Diffusion XL (SDXL) and other state-of-the-art text-to-image generation systems.

Image and Video Captioning: **Summary**

- ▶ Captioning lacks true understanding
- ▶ Multimodal alignment is shallow
- ▶ Hallucinations in T2I models
- ▶ Dataset bias and fairness issues
- ▶ High computational cost

- ▶ Vision-language reasoning (e.g., VQA with logic)
- ▶ 3D scene and video generation from text
- ▶ Interactive multimodal agents
- ▶ Multilingual and cultural alignment
- ▶ Personalized generation

Key Takeaways

- ▶ Vision and text integration unlocks powerful applications
- ▶ Techniques include captioning, CLIP, T2I, CFG
- ▶ Rapid evolution of generative multimodal models
- ▶ Ongoing challenges and exciting frontiers

Image and Video Captioning: **References**

- [1] Vinyals, O., Toshev, A., Bengio, S., & Erhan, D.
Show and Tell: A Neural Image Caption Generator.
In *CVPR*, 2015.
<https://arxiv.org/abs/1411.4555>
- [2] Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., & Sutskever, I.
Learning Transferable Visual Models From Natural Language Supervision (CLIP).
In *ICML*, 2021.
<https://arxiv.org/abs/2103.00020>

References (cont.)

- [3] Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., & Chen, M.
Hierarchical Text-Conditional Image Generation with CLIP Latents (DALL-E 2).
arXiv preprint, 2022.
<https://arxiv.org/abs/2204.06125>
- [4] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B.
High-Resolution Image Synthesis with Latent Diffusion Models (Stable Diffusion).
In *CVPR*, 2022.
<https://arxiv.org/abs/2112.10752>

References (cont.)

- [5] Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Salimans, T., Ho, J., Fleet, D. J., & Norouzi, M.
Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding (Imagen).
In *NeurIPS*, 2022.
<https://arxiv.org/abs/2205.11487>
- [6] BLIP: Bootstrapped Language-Image Pretraining.
<https://github.com/salesforce/BLIP>
- [7] OpenAI's DALL•E 2 Documentation.
<https://platform.openai.com/docs/guides/images>

References (cont.)

- [8] Midjourney Documentation.

<https://docs.midjourney.com/>

- [9] SDXL (Stable Diffusion XL) Documentation.

<https://stability.ai/guides/sdxl>

Credits

Dr. Prashant Aparajeya

Computer Vision Scientist — Director(AISimply Ltd)

p.aparajeya@aisimply.uk

This project benefited from external collaboration, and we acknowledge their contribution with gratitude.