# Vision-Language Pretraining: Current Trends and the Future

윤예준

# 목차

- Vision-Language Task
- Modern vision-language pretraining
- Current Trends and the Future

# Image Retrieval



"Grey haired man in black and yellow tie."

# Grounding Referring Expressions



"The man who is touching his head."

# Image Captioning



"A group of young people playing a game of Frisbee."

# Visual Question Answering(VQA)



Q: "What is the mustache made of?"  A: "bananas"

# Visual Dialog

# VL Datasets

- Image Retrieval: Flickr, COCO

- Grounding Referring Expression: RefCOCO, Visual7W

- Image Captioning: COCO

- Visual Question Answering: VQA v1, VQA v2, Visual Genome, GQA

- Visual Dialog: Visual Dialog, GuessWhat?!

# VL Datasets

- Image Retrieval: Flickr, COCO



① A child in a pink dress is climbing up a set of stairs in an entry way.

② A girl going into a wooden building.

③ A little girl climbing into a wooden playhouse.

④ A little girl climbing the stairs to her playhouse.

⑤ A little girl in a pink dress going into a wooden cabin.

# VL Datasets

• Grounding Referring Expression: RefCOCO, Visual7W



RefCOCO

woman on right in white shirt
woman on right
right woman

Q: What endangered animal is featured on the truck?
A: **A bald eagle.**
A: A sparrow.
A: A humming bird.
A: A raven.

Q: Where will the driver go if turning right?
A: **Onto 24 ¾ Rd.**
A: Onto 25 ¾ Rd.
A: Onto 23 ¾ Rd.
A: Onto Main Street.

Q: When was the picture taken?
A: **During a wedding.**
A: During a bar mitzvah.
A: During a funeral.
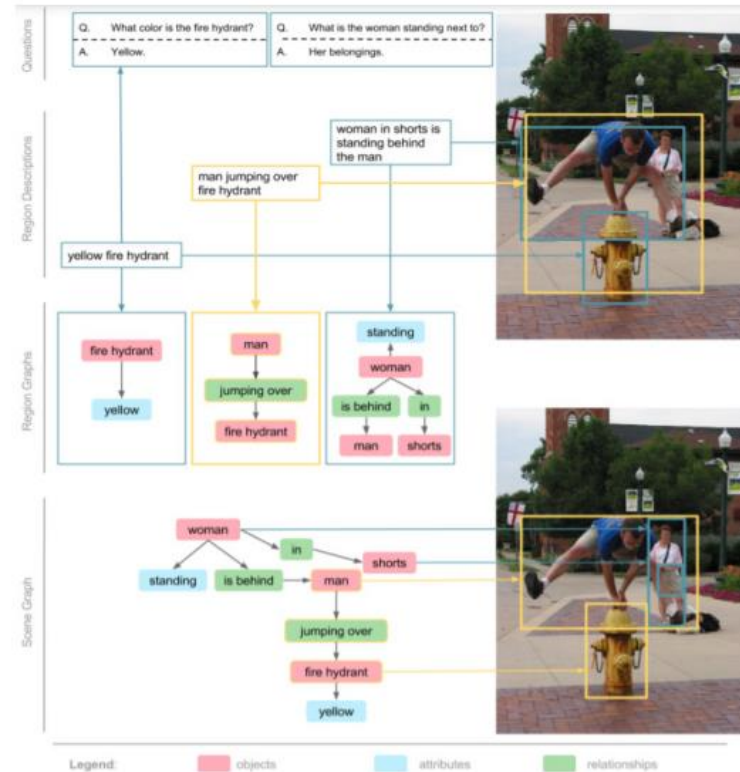A: During a Sunday church service.

Q: Who is under the umbrella?
A: **Two women.**
A: A child.
A: An old man.
A: A husband and a wife.

# VL Datasets

- Visual Question Answering: VQA v1, VQA v2, Visual Genome, GQA

# VL Datasets

- Image Captioning: [COCO](#)

# VL Datasets

- Visual Dialog: [Visual Dialog](#), [GuessWhat?!](#)



Caption: A man and woman on bicycles are looking at a map.
Person A (1): where are they located
Person B (1): in city
Person A (2): are they on road
Person B (2): sidewalk next to 1
Person A (3): any vehicles
Person B (3): 1 in background
Person A (4): any other people
Person B (4): no
Person A (5): what color bikes
Person B (5): 1 silver and 1 yellow
Person A (6): do they look old or new
Person B (6): new bikes
Person A (7): any buildings
Person B (7): yes
Person A (8): what color
Person B (8): brick
Person A (9): are they tall or short
Person B (9): i can't see enough of them to tell
Person A (10): do they look like couple
Person B (10): they are



**Questioner**
Is it a vase?
Is it partially visible?
Is it in the left corner?
Is it the turquoise and purple one?
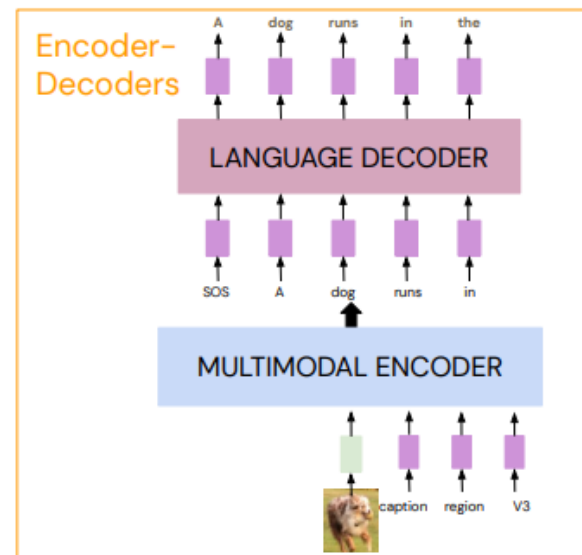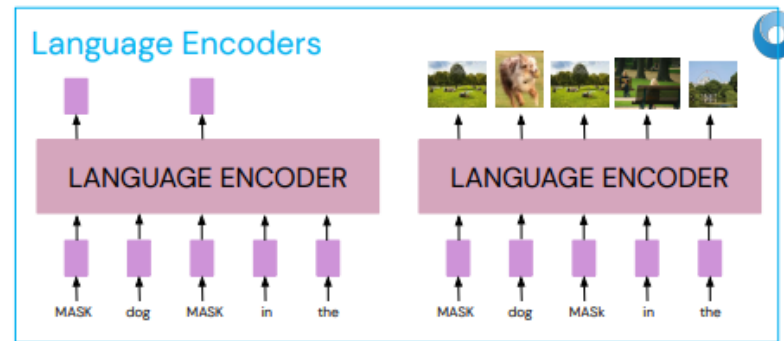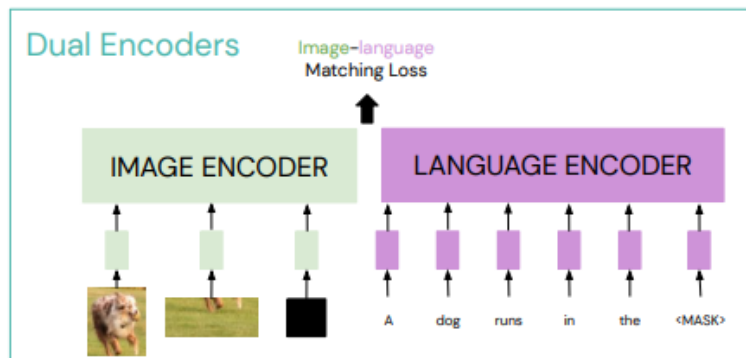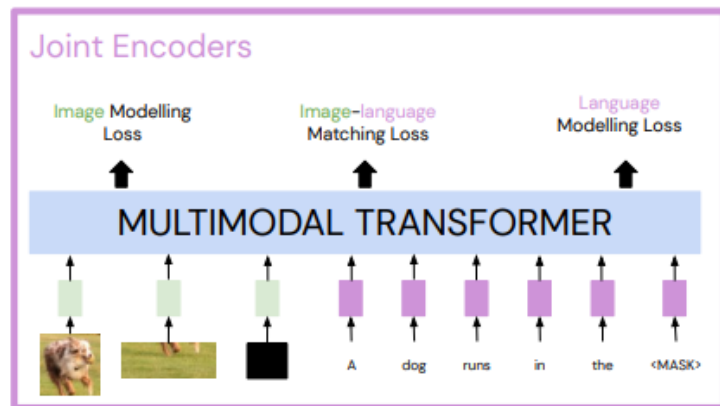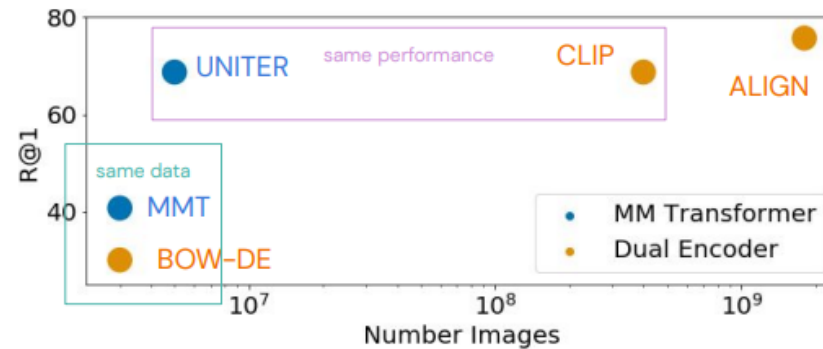
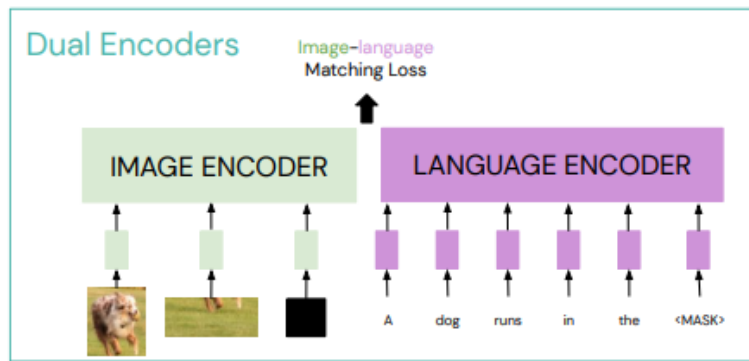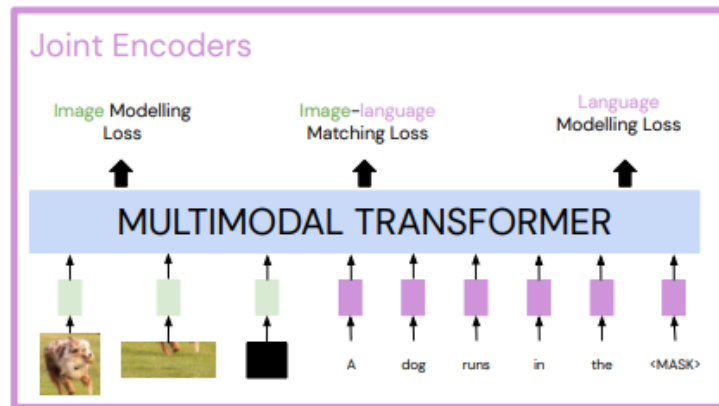**Oracle**
Yes
No
No
Yes

Figure 1: An example game. After a sequence of four questions, it becomes possible to locate the object (highlighted by a green bounding box).
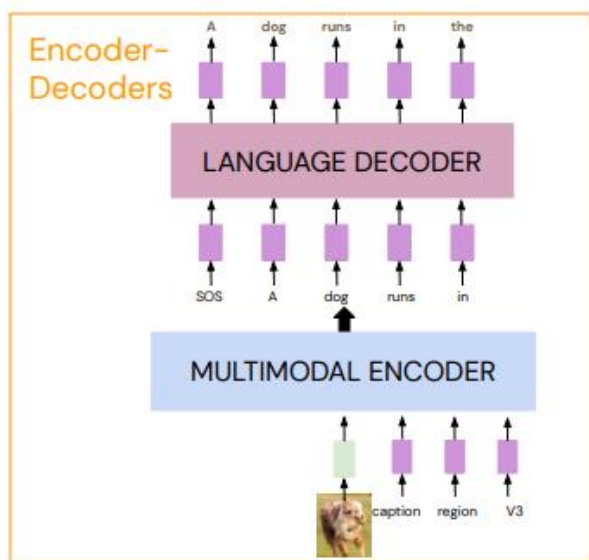
13

# 현대 모델 구조

# 현대 모델 구조

# 현대 모델 구조



**SimVLM** [Wang et al, 2022]

Unifies tasks as text generation.
Removes object detection supervision.
Trains on large-scale noisy image-text data (ALIGN).

# Current Trends and the Future

# Statistical learning has limitations.



**Predictions are reliable only within the training distribution.**

Challenging if the model relies on grass in the background.

Training data (biased)

Test data (out-of-distribution)

**The features used by a model are not necessarily the same as for the real system we try to imitate.** (e.g. human labeller)

Formally, in causal language: the background is not a *cause* to the image label.

⇔ Intervening on the background (by *editing the image*) would *not* cause one to label it differently.

A cow. → (Intervention) Still a cow.

# Causal learning

- Learning the data-generating mechanisms of a task (and not just the correlations in a specific dataset).

**Emerging area: extending ML with causal principles** (high-dimensional data & causal relationships not modelled explicitly)

› Causal representation learning: learning embeddings of raw data, disentangling its generative factors (causal parents).

*Equivalent to: disentanglement, independent component analysis (ICA).*

› Causal learning: learning predictive models that rely on causal (not spurious) features.

*Enable better transfer to unseen conditions, across datasets, across tasks.*

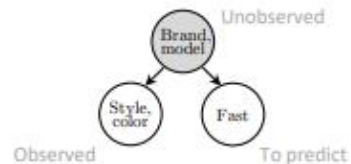*Also aims at (implicitly) identifying generative factors.*

# statistical vs. causal learning



**"Is this a fast car ?"** (top speed >200 km/h)

Training images with labels 'Fast' ∈ {0,1}

Someone's mental (causal) model:

Brand, model — Unobserved

Style, color — Observed

Fast — To predict

> **Statistical learning is about correlations:** red = fast.
> Reliable only if the training/test data are from similar distributions.

> **Causal learning is about mechanisms.**
> It enables predictions in conditions unobserved during training (OOD).

What happens to a re-painted car ? → **Faster ?** No !

Conditioned on **observing** the color in the training distribution.

$$\mathrm{P}(Fast \mid Color)$$

$$\neq$$

$$\mathrm{P}(Fast \mid \mathrm{do}(Color))$$

Conditioned on an **intervention**.

# statistical vs. causal learning

**Only 2 options to obtain knowledge of the data-generating process.**

> **Existing** task knowledge from humans.

Examples: custom architectures and losses,
hand-designed data augmentations,
interaction with human-designed simulator, etc.

> Heterogeneous/interventional data = **non-i.i.d. samples.**

Examples: data collected before/after interventions,
data from multiple environments (in time/location/subpopulation/...),
pairs of counterfactual examples,
non-stationary time series, etc.