# Multimodal Models: Part 1

Manish Gupta
Principal Applied Scientist, Microsoft
https://sites.google.com/view/manishg/

# Vision-and-Language Tasks



VQA

VCR Q→A          VCR QA→R

Referring Expressions

Caption-Based Image Retrieval

# Vision Transformers



**Vision Transformer (ViT)**

**Transformer Encoder**

Class: Bird, Ball, Car ...

MLP Head

Transformer Encoder

Patch + Position Embedding
* Extra learnable [class] embedding

Linear Projection of Flattened Patches

L ×
MLP
Norm
Multi-Head Attention
Norm
Embedded Patches

| Model | Layers | Hidden size $D$ | MLP size | Heads | Params |
|---|---|---|---|---|---|
| ViT-Base | 12 | 768 | 3072 | 12 | 86M |
| ViT-Large | 24 | 1024 | 4096 | 16 | 307M |
| ViT-Huge | 32 | 1280 | 5120 | 16 | 632M |

- Split an image into fixed-size patches
- Linearly embed each of them
- Add 1D position embeddings
- Feed the resulting sequence of vectors (prepended by [CLS]) to a standard Transformer encoder.
- Classification MLP head with 1 hidden layer at pre-training and just a single linear layer at fine-tune time.
- Pretrain datasets: ImageNet-1K, ImageNet-21k, JFT
- ViT-L/16 means the "Large" variant with 16×16 input patch size. Smaller path size ➜ larger seq length.
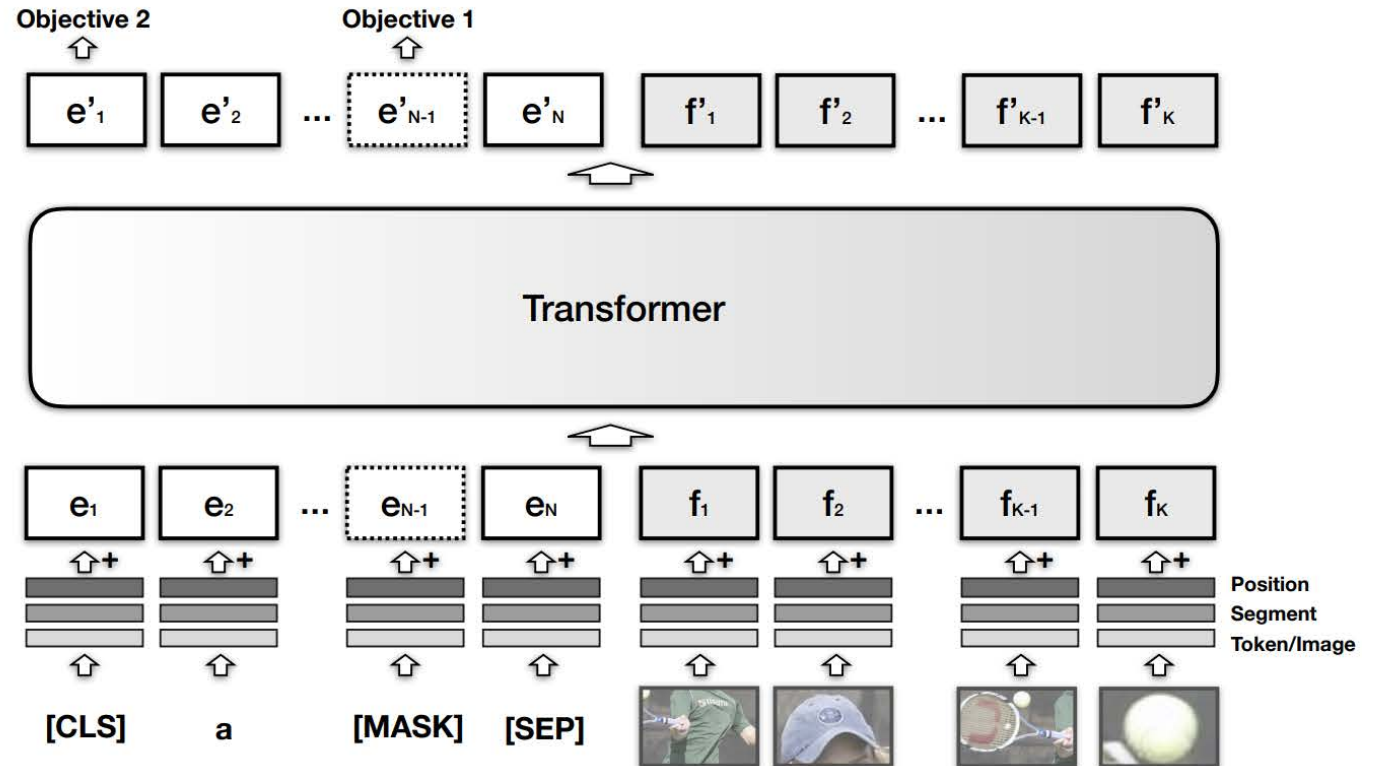- Match or exceed accuracy of ResNets on many image classification datasets

Dosovitskiy, Alexey, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani et al. "An image is worth 16x16 words: Transformers for image recognition at scale." arXiv:2010.11929 (2020).

# Joint representation model for vision and language
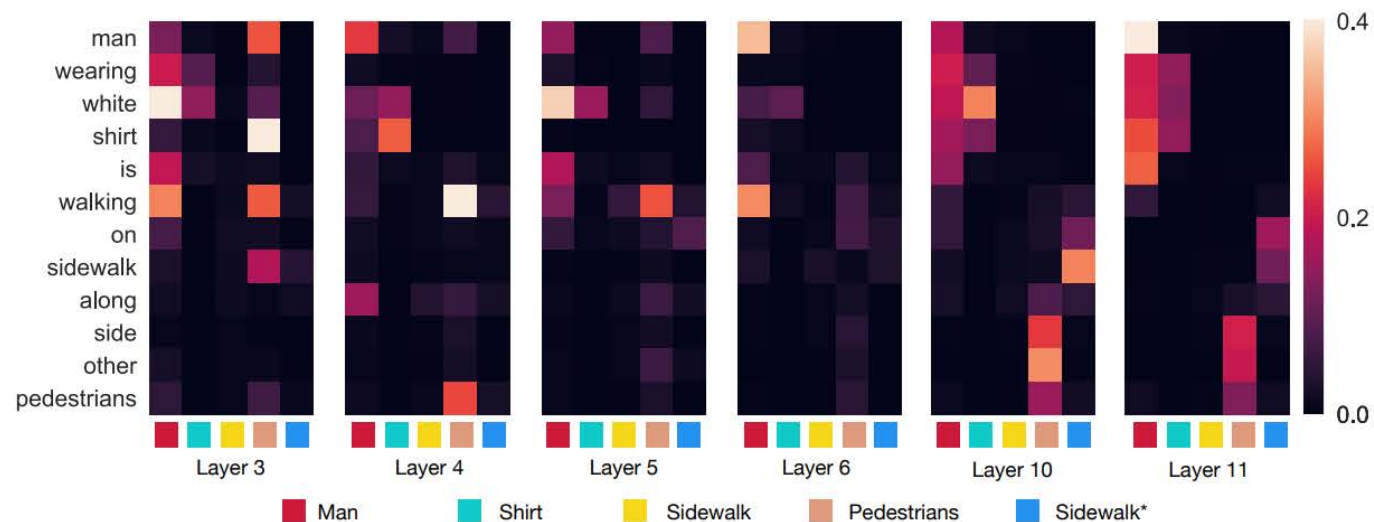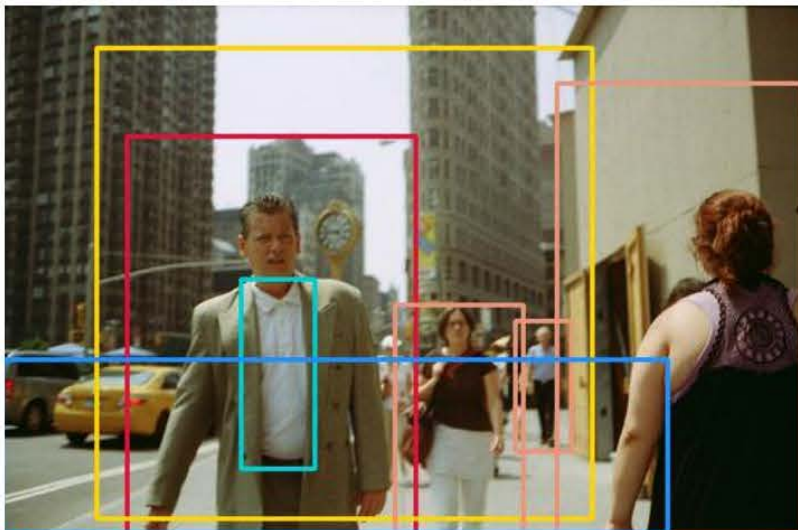


A person hits a ball with a tennis racket

- MLM (Objective 1), and sentence-image prediction task (Objective 2)
- VisualBERT integrates BERT for NLP, and pretrained object proposals systems such as Faster-RCNN.

Li, Liunian Harold, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. "Visualbert: A simple and performant baseline for vision and language." arXiv:1908.03557 (2019).
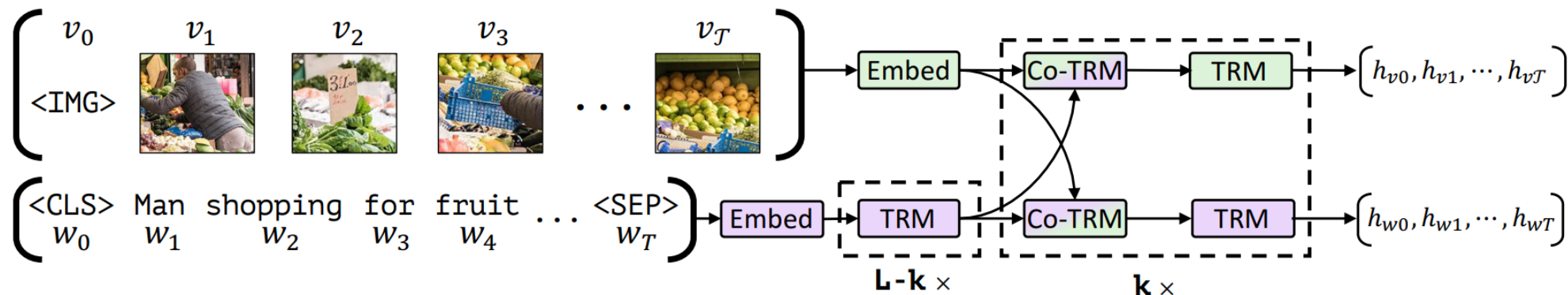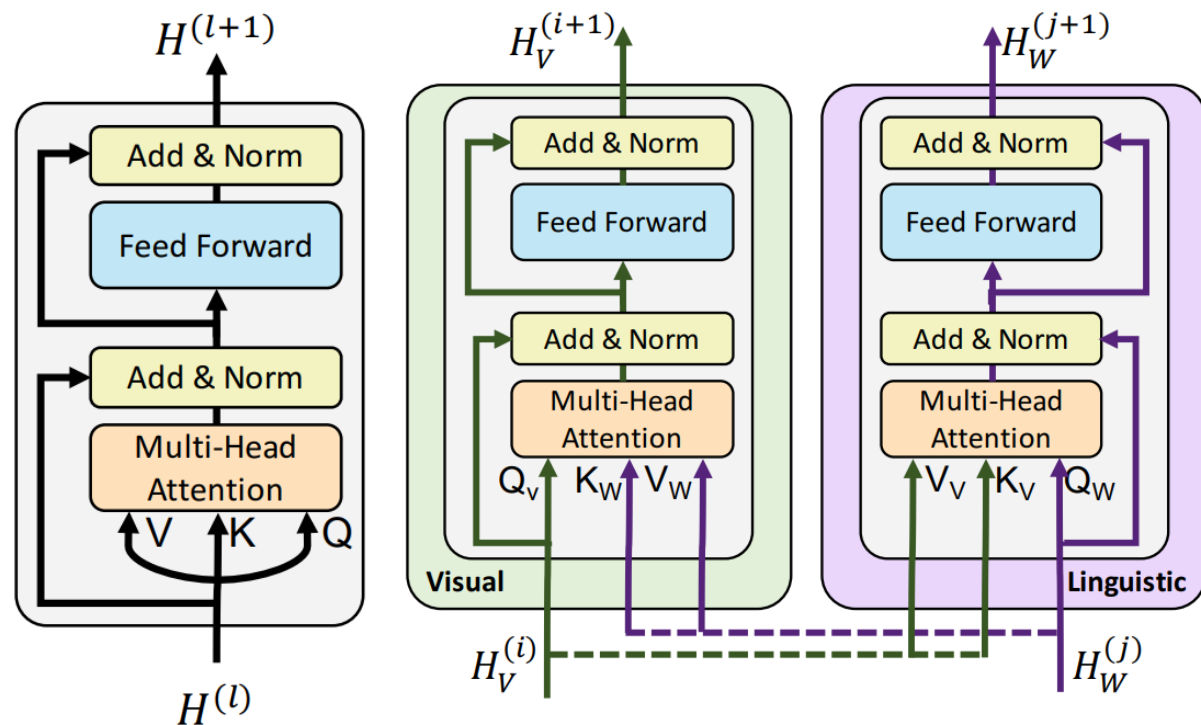
Manish Gupta

# VisualBERT



- Attention weights of some selected heads in VisualBERT.
- In high layers (e.g., the 10-th and 11-th layer), VisualBERT is capable of implicitly grounding visual concepts (e.g., "other pedestrians" and "man wearing white shirt").
- The model also refines its understanding over the layers, incorrectly aligning "man" and "shirt" in the 3-rd layer but correcting them in higher layers.

Li, Liunian Harold, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. "Visualbert: A simple and performant baseline for vision and language." arXiv:1908.03557 (2019).

LLMs: Introduction and Recent Advances                    Manish Gupta
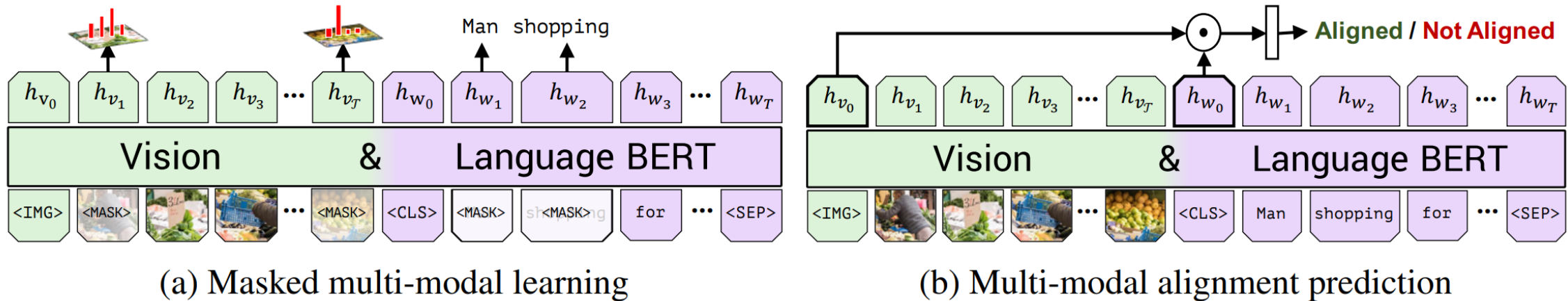
# ViLBERT Architecture

- The text stream has significantly more processing before interacting with visual features.

- Initialize the linguistic stream of ViLBERT with BERT BASE. Use Faster R-CNN pretrained on the Visual Genome dataset.



Lu, Jiasen, Dhruv Batra, Devi Parikh, and Stefan Lee. "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks." arXiv:1908.02265 (2019).

Manish Gupta

# ViLBERT Training Tasks and Objectives



(a) Masked multi-modal learning
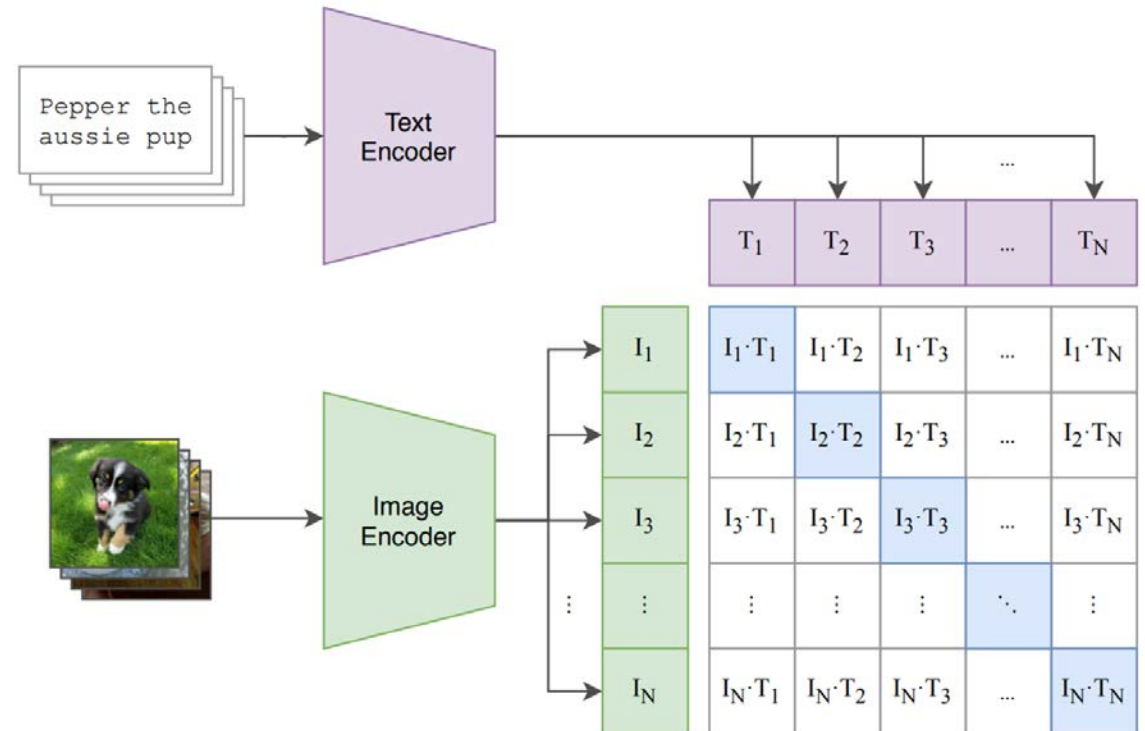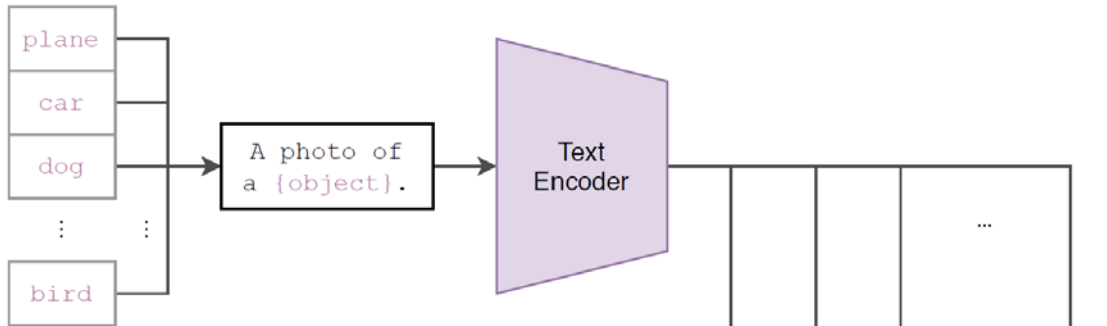
(b) Multi-modal alignment prediction

- Train ViLBERT on Conceptual Captions (~3.3M images) to learn visual grounding.
- Masked multi-modal learning: reconstruct image region categories or words for masked inputs given the observed inputs.
- Multi-modal alignment prediction: predict whether or not the caption describes the image content.

Lu, Jiasen, Dhruv Batra, Devi Parikh, and Stefan Lee. "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks." arXiv:1908.02265 (2019).

# CLIP (Contrastive Language-Image Pre-training)

- Pre-trained using WebImageText (WIT) 400M (image, text) pairs.

- Text encoder is a 12L Transformer.

- 5 ResNets
  - ResNet-50, a ResNet-101
  - RN50x4, RN50x16, and RN50x64: use ~4x, 16x, and 64x the compute of a ResNet-50.

- 3 Vision Transformers (ViT)
  - ViT-B/32, a ViT-B/16, and a ViT-L/14

- Maximize cos-sim of the image and text embeddings of N real pairs in the batch

- Minimize cos-sim of the embeddings of the N × N – N incorrect pairings.

- Tested on 30+ CV tasks like OCR, action recognition in videos, geo-localization,...

- 0-shot CLIP is often ≡ fully supervised baseline

# CLIP (Contrastive Language-Image Pre-training)



(2) Create dataset classifier from label text

(3) Use for zero-shot prediction

A zero-shot CLIP classifier outperforms a fully supervised linear classifier fitted on ResNet-50 features on 16 datasets, including ImageNet.

# Classification using CLIP



**FOOD101**

**guacamole** (90.1%) Ranked 1 out of 101 labels

✓ a photo of **guacamole**, a type of food.

✗ a photo of **ceviche**, a type of food.

✗ a photo of **edamame**, a type of food.

✗ a photo of **tuna tartare**, a type of food.

✗ a photo of **hummus**, a type of food.

**SUN397**

**television studio** (90.2%) Ranked 1 out of 397

✓ a photo of a **television studio**.

✗ a photo of a **podium indoor**.

✗ a photo of a **conference room**.

✗ a photo of a **lecture room**.

✗ a photo of a **control room**.

**YOUTUBE-BB**

**airplane, person** (89.0%) Ranked 1 out of 23

✓ a photo of a **airplane**.

✗ a photo of a **bird**.

✗ a photo of a **bear**.

✗ a photo of a **giraffe**.

✗ a photo of a **car**.

**EUROSAT**

**annual crop land** (12.9%) Ranked 4 out of 10

✗ a centered satellite photo of **permanent crop land**.

✗ a centered satellite photo of **pasture land**.

✗ a centered satellite photo of **highway or road**.

✓ a centered satellite photo of **annual crop land**.

✗ a centered satellite photo of **brushland or shrubland**.

# Visually-rich Document Understanding



Q: Mention the ZIP code written?
A: 80202
Q: What date is seen on the seal at the top of the letter?
A: 23 sep 1970
Q: Which company address is mentioned on the letter?
A: Great western sugar Co.

Xu, Yang, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu et al. "LayoutLMv2: Multi-modal Pre-training for Visually-rich Document Understanding." ACL, pp. 2579-2591. 2021.

# Visually-rich Document Understanding

- Xu, Yang, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Y... "LayoutLMv2: Multi-modal Pre-training for Visually-rich Document U... pp. 2579-2591. 2021.

Manish Gupta

# LayoutLMv2 Architecture

- Text: Initialized using UniLMv2

- ResNeXt-FPN architecture with MaskRCNN backbone of the visual encoder.

- Use output feature map (W=H=7).

- Embed spatial layout of token bounding boxes from the OCR results
  - Concat(PosEmb2Dx($x_{min}$, $x_{max}$, width), PosEmb2Dy($y_{min}$, $y_{max}$, height))

- LayoutLMv2
  - Base: 12 layers (200M)
  - Large: 24 layers (426M)

- 3 tasks: MVLM, TIA, TIM

- Dataset: 11M scanned docs. Text OCR: Microsoft Read API



Xu, Yang, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu et al. "LayoutLMv2: Multi-modal Pre-training for Visually-rich Document Understanding." ACL, pp. 2579-2591. 2021.

# Video tasks

- Text→Video Retrieval
  - Given text and a collection of videos, find relevant ones.

- Multiple-choice VideoQA.
  - Given video, a question and multiple candidate answers, choose the best one.

- Action Segmentation/Action Step Localization
  - Assign each token (or frame) of a video with one of the pre-defined labels (or steps) to separate meaningful segments of videos.
  - Similar to sequence labeling (e.g. NER) in NLP.
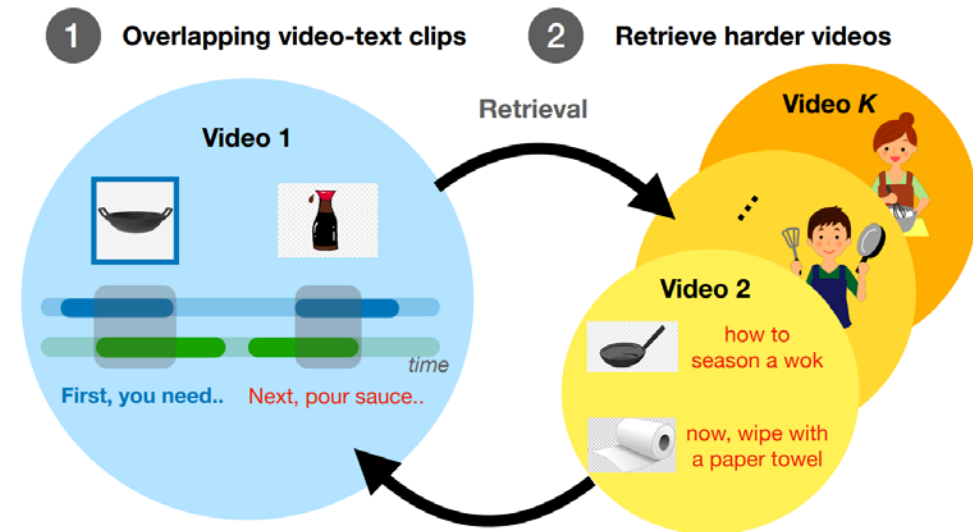
Xu, Hu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer.

"VideoCLIP: Contrastive Pre-training for Zero-shot Video-Text Understanding." In EMNLP, pp. 6787-6800. 2021.

# What is the VideoCLIP architecture? How it is pretrained?

- Contrastive approach to pre-train a unified model for zeroshot video and text understanding.
  - Loosely temporally overlapping positive video-text pairs, instead of enforcing strict start/end timestamp overlap.
  - Hard negatives using nearest neighbor retrieval that uses video clusters to form batches with mutually harder videos.
- BERT-base-uncased for both video (6L) and text (12L).
- Video: frozen pretrained CNN, projected to video tokens using a MLP layer.
- Average pooling over the sequence of tokens for video and text.
- Pretraining data: HowTo100M



VideoCLIP: Contrastive learning with hard-retrieved negatives and overlapping positives for video-text pre-training.

$$\mathcal{L} = -\sum_{(v,t)\in B} \Big( \log \text{NCE}(z_v, z_t) + \log \text{NCE}(z_t, z_v) \Big)$$

$$\text{NCE}(z_v, z_t) = \frac{\exp\left(z_v \cdot z_t^+/\tau\right)}{\sum_{z\in\{z_t^+, z_t^-\}} \exp\left(z_v \cdot z/\tau\right)}$$

Xu, Hu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. "VideoCLIP: Contrastive Pre-training for Zero-shot Video-Text Understanding." In EMNLP, pp. 6787-6800. 2021.

# What is ImageBind?

- An image of a beach can remind us of the sound of waves, the texture of the sand, a breeze, or even inspire a poem.

- Aligns six modalities' embedding into a common space: images, text, audio, depth, thermal, and Inertial Measurement Unit (IMU).

- Image-paired data is sufficient to bind the modalities together.



Girdhar, Rohit, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. "Imagebind: One embedding space to bind them all." In CVPR, pp. 15180-15190. 2023.

# How is the ImageBind model trained?

$$L_{\mathcal{I},\mathcal{M}} = -\log \frac{\exp(\mathbf{q}_i^\mathsf{T}\mathbf{k}_i/\tau)}{\exp(\mathbf{q}_i^\mathsf{T}\mathbf{k}_i/\tau) + \sum_{j\neq i}\exp(\mathbf{q}_i^\mathsf{T}\mathbf{k}_j/\tau)}$$



- $q_i = f(I_i)$ and $k_i = g(M_i)$ where f and g are deep networks.
- InfoNCE loss. Symmetric loss $L_{I,M} + L_{M,I}$
- ViT-H 630M params; text encoders (302M params) from OpenCLIP (frozen)
- Same encoder for images+videos. Treat videos as multi-frame images.

- Datasets:
  - (video, audio) pairs from Audioset
  - (image, depth) pairs from SUN RGB-D
  - (image, thermal) pairs from LLVIP
  - (video, IMU) pairs from Ego4D
  - (image, text) pairs from large-scale web data.

Girdhar, Rohit, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. "Imagebind: One embedding space to bind them all." In CVPR, pp. 15180-15190. 2023.

# Thanks!

- HomePage: https://sites.google.com/view/manishg/

- Google Scholar: https://scholar.google.co.in/citations?user=eX9PSu0AAAAJ

- LinkedIn: http://aka.ms/manishgupta

- YouTube (Data Science Gems): https://www.youtube.com/@dlByManish