# "Lecture" 13:
# Course Conclusion

# Announcements

- **Please make sure to sign up for project advising sessions today and tomorrow** (see Ed post for the sign-up sheet), this is 5% of your project grade. Contact Xiaohan if an alternate time is needed.
- Project milestone grades will be provided before the session time.
- **Project final presentation is during the final exam slot Mon Dec 9** from 3:30-6:30pm. More details will be provided, expect ~7 minutes presentation per project.
- **Project final report is due Wed Dec 11**

- James's office hours will start at 5:30 today instead of 4pm; if this poses difficulty for you, please contact him to arrange an alternate time
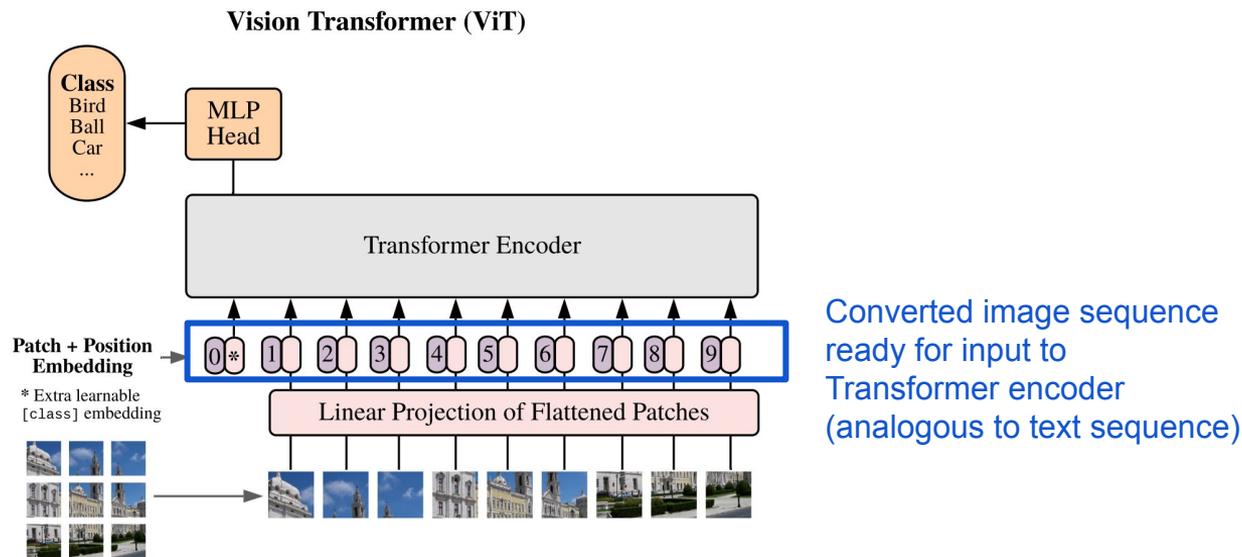
# Announcements

- **Wednesday will be a final guest lecture from Khaled Saab, first-author of the Med-Gemini paper. This lecture will be remote over zoom (link will be provided in Ed), please plan to attend!**

# Our goals at the beginning of the course

- Fluency in cutting edge computer vision models and research
  - Focus on vision and vision-language foundation models, from representation learners to diffusion and generative models
- Understanding of existing applications and opportunities for future utility in biomedicine
- Ability to think through considerations for real-world use, including model size and computation, training and inference settings, and training data

# First topic: Vision Transformers (ViT)

Key idea: Convert image into sequence of patches. Can then benefit from Transformer architecture and self-attention, which jointly attends over all patches



Dosovitsky et al. 2021

# Two major classes of vision foundation models

## Representation Learners

Train on huge amounts of data

Objective: learn powerful generalist feature extractors that can be used for downstream computer vision tasks

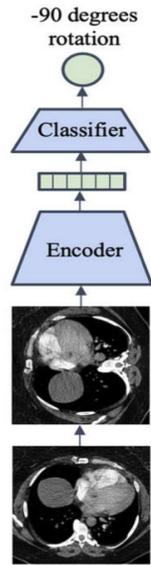- Vision Only (e.g. DINOv2)

- Vision-Language (e.g. CLIP)

## Generative Models

Train on huge amounts of data

Objective: learn to generate images and/or text

- Text -> Image (e.g. StableDiffusion, DALL-E)

- Text, Image -> Text (e.g. Flamingo, GPT4-V)

# Two major classes of vision foundation models

## Representation Learners

Train on huge amounts of data

Objective: learn powerful generalist feature extractors that can be used for downstream computer vision tasks

- Vision Only (e.g. DINOv2)

- Vision-Language (e.g. CLIP)

## Generative Models
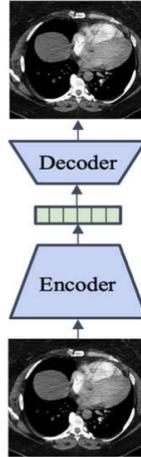
Train on huge amounts of data

Objective: learn to generate images and/or text

- Text -> Image (e.g. StableDiffusion, DALL-E)

- Text, Image -> Text (e.g. Flamingo, GPT4-V)
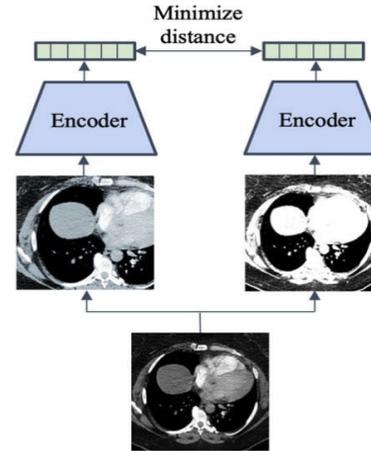
# Different representation learning paradigms



**Popular state-of-the-art approaches**

-90 degrees rotation

Minimize distance

**Innate relationship objective**
E.g., predict rotation angle (or some other innate property) of an image

**Generative objective**
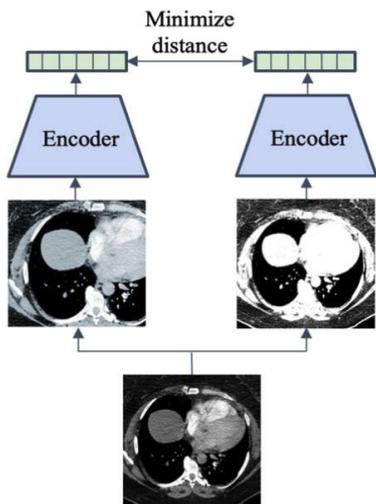Compress and then reconstruct input image (e.g. autoencoders)

**Contrastive objective**
Different views of the same input should have more similar representation to each other than with a different input

**Self-prediction objective**
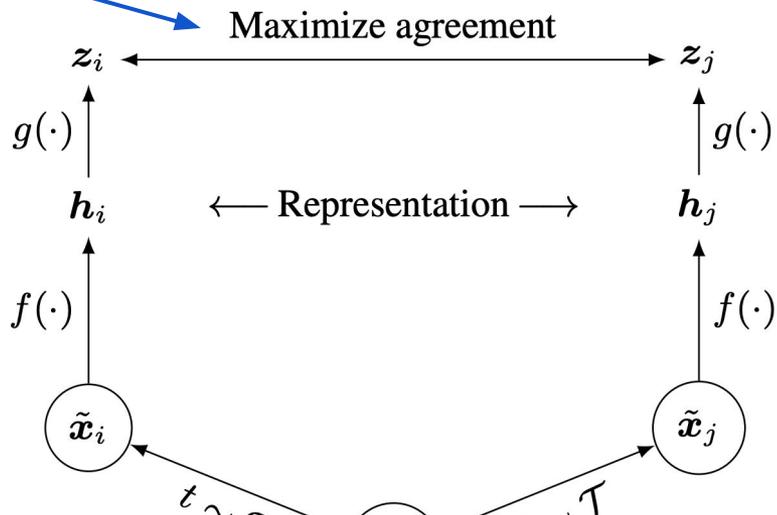Mask parts of input data and predict these parts

Figure credit: Huang et al. 2023.

# SimCLR: a foundational method for contrastive self-supervised learning



Minimize distance

Encoder    Encoder

Contrastive loss

**SimCLR formulation**

Maximize agreement

$z_i \longleftrightarrow z_j$

$g(\cdot) \uparrow \qquad \uparrow g(\cdot)$

$h_i \quad \longleftarrow \text{Representation} \longrightarrow \quad h_j$

$f(\cdot) \uparrow \qquad \uparrow f(\cdot)$

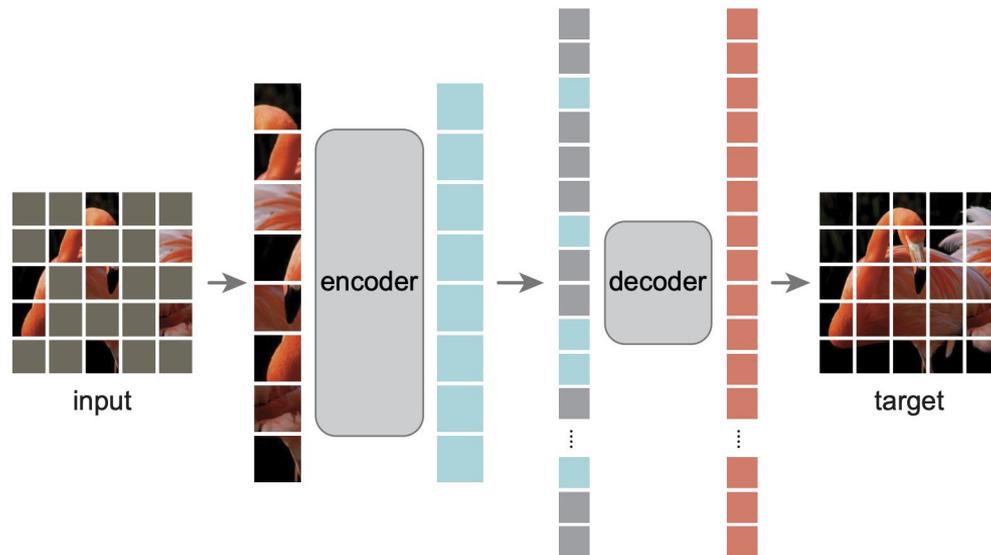$\tilde{x}_i \qquad\qquad \tilde{x}_j$

$t \sim \mathcal{T} \qquad \qquad \mathcal{T}$

After self-supervised training, can fine-tune the encoder *f* on smaller labeled datasets. Can also directly extract learned representations h for downstream tasks.
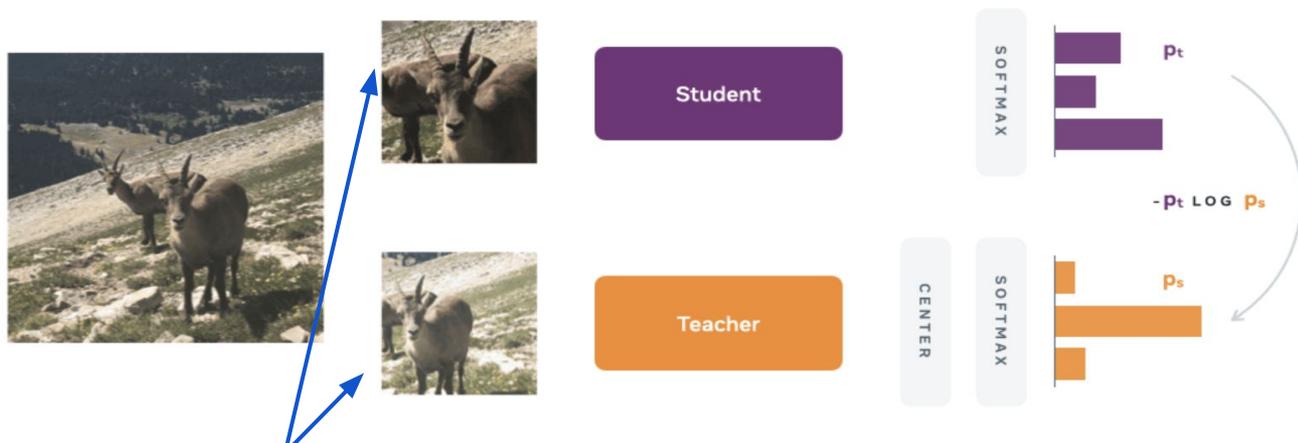
# Masked Autoencoders (MAE)

- Key idea: mask substantial parts of the input, train the model to reconstruct (predict) these parts

- Inspired by major self-supervised representation learning paradigm in NLP (e.g. BERT), that masks tokens in sentences and trains models to reconstruct them



He et al. 2021

# Another SSL approach: DINO (Self-**Di**stillation with **No** Labels)

Not an example of a contrastive learning objective for self-supervised learning, but related! Builds on the notion of matching representations of augmented views of the same image, but no longer uses negative samples. Instead, proposes a **teacher-student framework**.
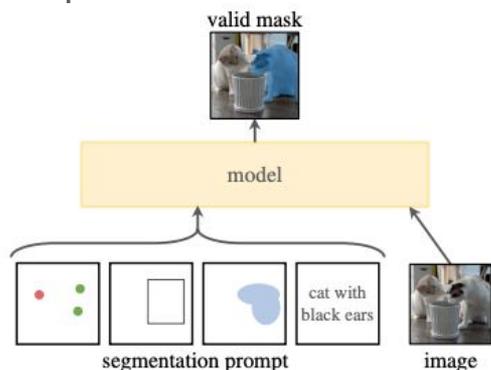


Augmented "views" of the input are passed to both the student and teacher networks.
**Student views**: more diverse and aggressive augmentations, to learn representations that generalize across augmentations
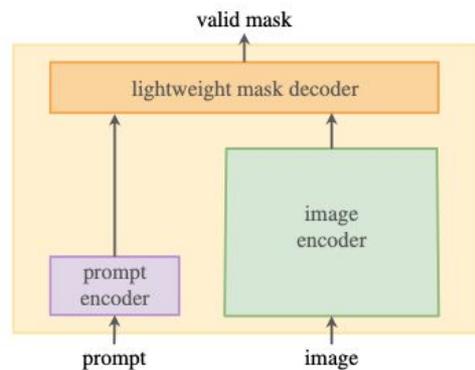**Teacher views**: less aggressive augmentations (including larger crops), to provide a more stable target for the student to follow.

Caron et al. 2021
Oquab et al. 2024

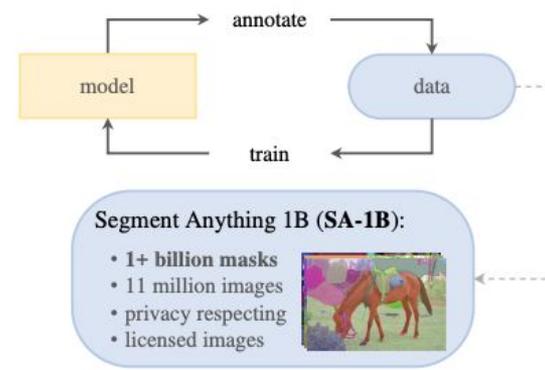# Segment Anything Model (SAM): A foundation model targeted for segmentation

- Foundation model for promptable segmentation, based on a Transformer encoder-decoder architecture. Generalizes to many segmentation tasks.

- Not all representation learning needs to be through self-supervised learning. Here, a supervised paradigm that also achieves powerful representation learning.

- Trained on 1 billion masks from 11 million images, using the model in a data collection loop.



Kirrilov et al. 2023.

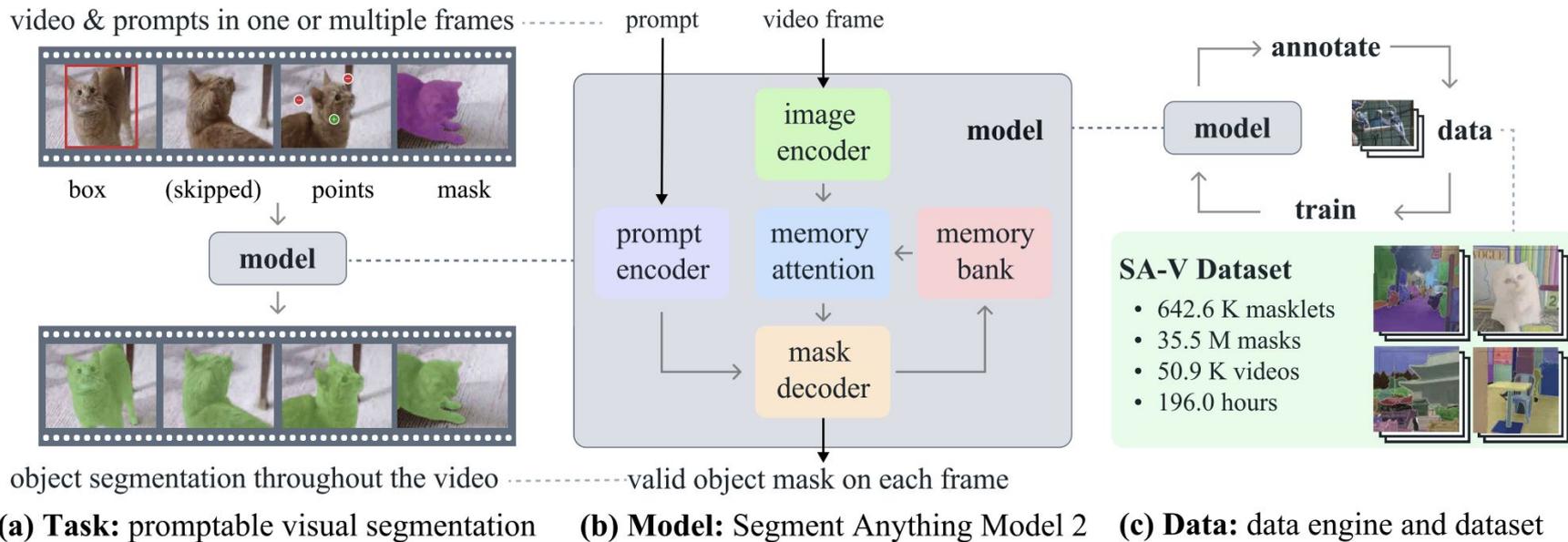(a) **Task**: promptable segmentation     (b) **Model**: Segment Anything Model (**SAM**)     (c) **Data**: data engine (top) & dataset (bottom)

# Segment Anything Model 2 (SAM 2): Extending to videos



**(a) Task:** promptable visual segmentation

**(b) Model:** Segment Anything Model 2

**(c) Data:** data engine and dataset

SA-V Dataset
- 642.6 K masklets
- 35.5 M masks
- 50.9 K videos
- 196.0 hours

Ravi et al. 2024

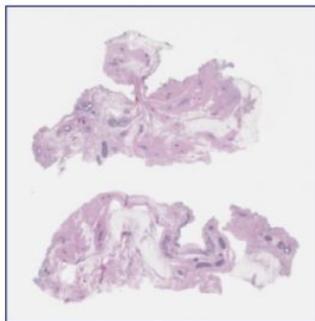# RAD-DINO: Improving scale and model for chest X-ray pre-training

- Used the DINOv2 self-supervised method, with ViT architecture. Started from public DINOv2 ViT-B/14 model and continued training with CXR images.

- Trained on an extended CXR dataset, Multi-CXR, aggregating multiple datasets including CheXpert

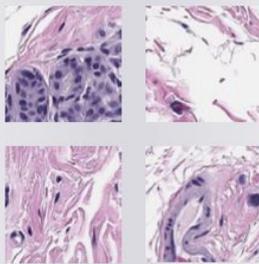| Dataset | View | Patient cohort | Number of subjects | Number of images |
|---|---|---|---|---|
| BRAX [122] | frontal, lateral | all available in institutional PACS | 19,351 | 41,620 |
| CheXpert [69] | frontal, lateral | inpatient and outpatient | 65,240 | 223,648 |
| MIMIC-CXR [17] | frontal | ICU | 188,546 | 210,491 |
| NIH-CXR [111] | frontal | not specified | 32,717 | 112,120 |
| PadChest [57] | frontal, lateral | all available | 67,000 | 160,817 |
| Private | frontal, lateral | outpatient | 66,323 | 90,000 |
| **Total** | | | 439,177 | 838,336 |

Perez-Garcia et al. RAD-DINO: Exploring scalable medical image encoders beyond text supervision. 2024.

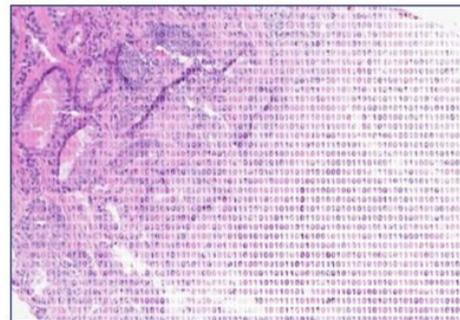# Virchow: Large-scale self-supervised learning on digital pathology images

- Trained a DINO v2 based model on 1.5 million H&E stained whole slide images (WSIs) from Memorial Sloan Kettering Cancer Center and external consults, corresponding to ~100,000 patients



H&E slide

Tissue tiles
224 × 224 pixel crops crops from
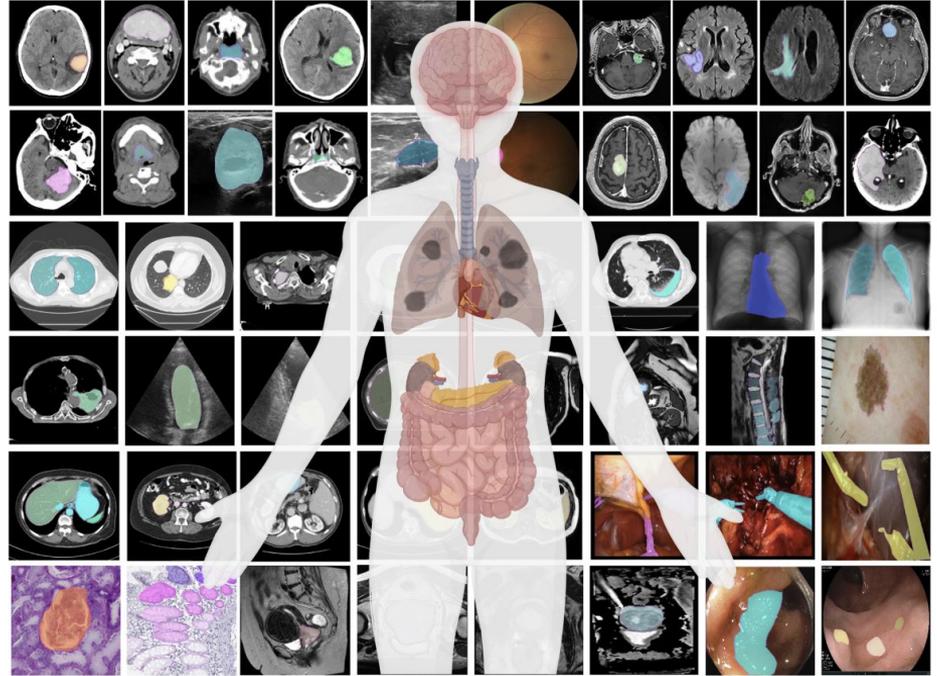tissue regions in the slide

Virchow
Foundation model with ViT-H architecture (632 million
parameters) trained using DINOv2 framework

Vorontsov et al. A foundation model for clinical-grade computational pathology and rare cancers detection. Nature Medicine, 2024.

# MedSAM: Segment anything in medical images

- Tackles universal medical image segmentation across many imaging modalities and tasks

- Trained on a large-scale medical image dataset with 1,570,263 image-mask pairs, covering 10 imaging modalities and over 30 cancer types

- Evaluated on 86 internal validation and 60 external validation tasks



Ma et al. Segment anything in medical images. Nature Communications, 2024.

# Two major classes of vision foundation models

## Representation Learners

Train on huge amounts of data

Objective: learn powerful generalist feature extractors that can be used for downstream computer vision tasks

- Vision Only (e.g. DINOv2)

- Vision-Language (e.g. CLIP)

## Generative Models

Train on huge amounts of data

Objective: learn to generate images and/or text

- Text -> Image (e.g. StableDiffusion, DALL-E)

- Text, Image -> Text (e.g. Flamingo, GPT4-V)

# Two major classes of vision foundation models

**Representation Learners**

Train on huge amounts of data

Objective: learn powerful generalist feature extractors that can be used for downstream computer vision tasks

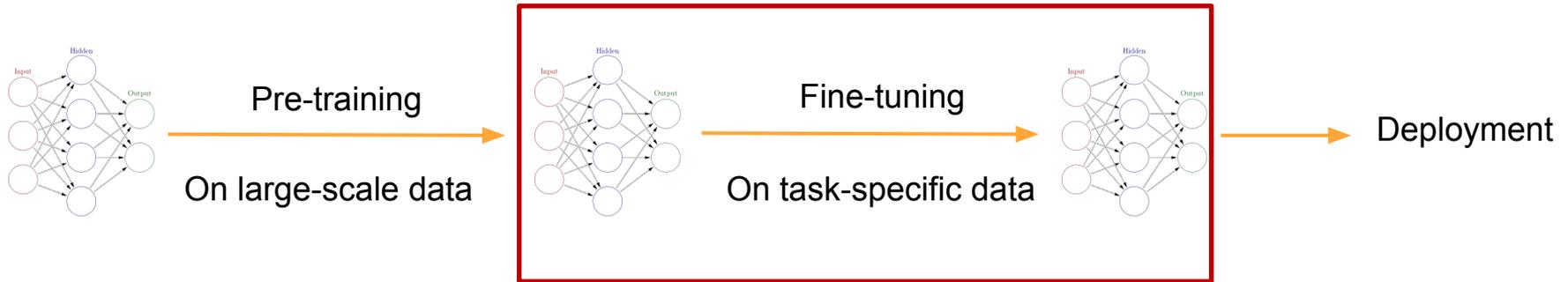- Vision Only (e.g. DINOv2)

- Vision-Language (e.g. CLIP)

**Generative Models**

Train on huge amounts of data
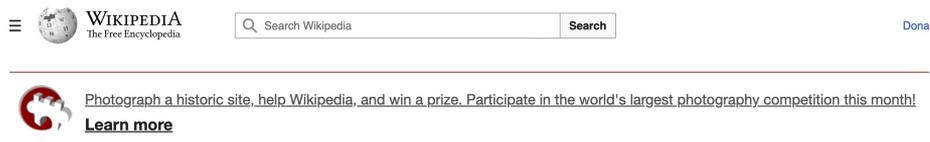
Objective: learn to generate images and/or text

- Text -> Image (e.g. StableDiffusion, DALL-E)

- Text, Image -> Text (e.g. Flamingo, GPT4-V)

# Adapt to Downstream Task



Pre-training

On large-scale data

Fine-tuning

On task-specific data

Deployment

How to skip this step to enable zero-shot generalization?
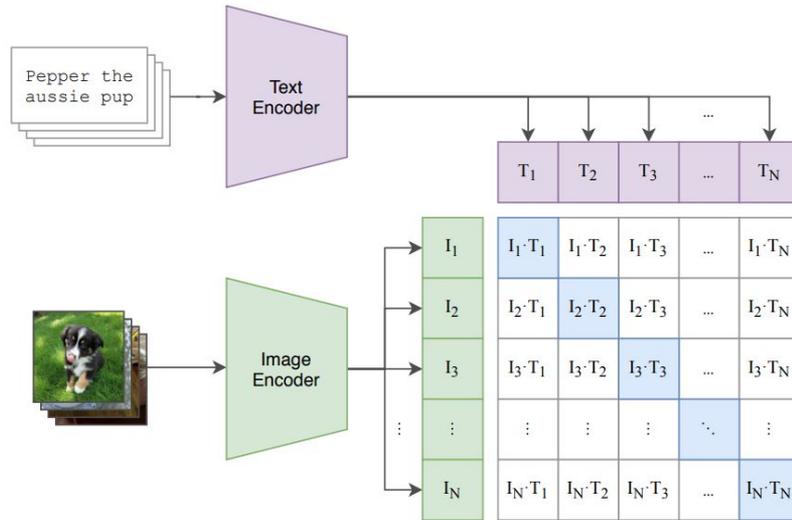
# Natural Language Supervision



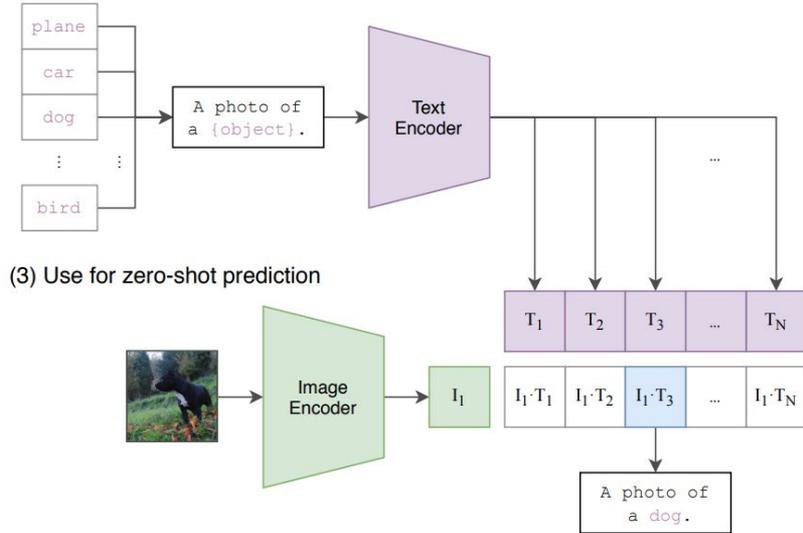❖ Unlimited Data: vast amount of text-image pairs on the internet

❖ Flexible zero-shot transfer: learn a representation that connects visual content to language

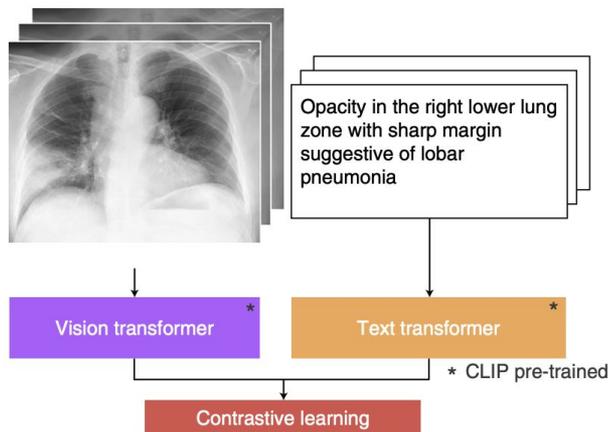# CLIP learns a joint representation space for images and text



Radford et al. 2021.

# CheXzero: leveraged CLIP to improve on ConVIRT and demonstrate zero-shot capabilities

- Compared to ConVIRT, updated to the same architecture as CLIP (better encoders, no nonlinear projection) as well as CLIP-pretrained weights.
- Also used the "impressions" section of the radiology report instead of ConVIRT sampling.



Tiu et al. Expert-level detection of pathologies from unannotated chest X-ray images via self-supervised learning. Nature Biomedical Engineering, 2022.

# PLIP: vision-language foundation model for pathology trained from Twitter

- Curated and trained on OpenPath: 208,414 image-text pairs scraped from the Internet, mostly from Twitter



Huang et al. A visual–language foundation model for pathology image analysis using medical Twitter. Nature Medicine, 2023.

# MONET: Leveraging a contrastively trained model to perform dataset and model auditing



Kim et al. Transparent medical image AI via an image–text foundation model grounded in medical literature. Nature Medicine, 2024.

# BiomedCLIP: Generalist foundation model trained on PubMed

PMC-15M: 15 million image-caption pairs from 4.4 million publicly available full-text articles in PubMed Central

Complete article packages are downloaded, and figure files and matching captions are extracted



Zhang et al. BiomedCLIP: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs. arXiv, 2024.

# Two major classes of vision foundation models

## Representation Learners

Train on huge amounts of data

Objective: learn powerful generalist feature extractors that can be used for downstream computer vision tasks

- Vision Only (e.g. DINOv2)

- Vision-Language (e.g. CLIP)

## Generative Models

Train on huge amounts of data

Objective: learn to generate images and/or text

- Text -> Image (e.g. StableDiffusion, DALL-E)

- Text, Image -> Text (e.g. Flamingo, GPT4-V)

# Diffusion training as learning to reverse noising

$$p(z_{t-1}|z_t)$$

$z_T \rightarrow \cdots \rightarrow z_t \rightarrow z_{t-1} \rightarrow \cdots \rightarrow x$

$$q(z_t|z_{t-1})$$

Sample t,
Apply noise: $q(z_t|z_{t-1})$



$z_{t-1}$  +  =  $z_t$

Use Unet to reverse the noise



$z_t$

$t$

$\hat{z}_\theta(z_t, t)$

Reconstruction loss

$$\mathbb{E}_{t \sim U(2,T)} \left[ c(\alpha_t) \| \hat{z}_\theta(z_t, t) - z_{t-1} \| \right]$$

# DDPM: the original pixel-level diffusion model



**Denoising Diffusion Probabilistic Models**

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu

Pieter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Figure 2: The directed graphical model considered in this work.

# Stable Diffusion or Latent Variable Models

Idea: instead of diffusion models on *pixels*, do diffusion in a *compressed latent space*



Diffusion
(DDPM)

Latent
diffusion
(Stable
Diffusion)

"High-Resolution Image Synthesis with Latent Diffusion Models"
"Score-based Generative Modeling in Latent Space"

# Stable Diffusion: example of text -> image generative model (+ other types of conditioning!)

Text -> image model trained on LAION: open-source dataset of 400M image-text pairs

Rombach et al. 2022.

# DALL-E: related family of OpenAI models that focus on high quality text-conditioned image generation

DALL-E 3 is trained on large amounts of detailed, synthetically generated captions to enable high-fidelity and high-detail generation



https://openai.com/dall-e-3

# One widely considered application area of vision generative models: augmenting training data

- RoentGen fine-tunes the Stable Diffusion model on the MIMIC-CXR dataset of chest x-ray (CXR) images and corresponding text reports (about ~175k images for training)



Bluethgen et al. A vision–language foundation model for the generation of realistic chest X-ray images. Nature Biomedical Engineering 2024.

# Preview to discussion paper presentations: generating pathology whole-slide image titles conditioned on RNA-sequencing



Carrillo-Perez et al. Generation of synthetic whole-slide image tiles of tumours from RNA-sequencing data via cascaded diffusion models. Nature Biomedical Engineering 2024.

# Diffusion models for MRI denoising: DDM$^2$

- MRI denoising is a key problem in modern MRI research due to the trade-off between achieving high signal-to-noise MRI scans and requiring long scan time (increased cost and discomfort, harder to accommodate overall patient demand)
- DDM$^2$ (Denoising Diffusion Models for Denoising Diffusion MRI) uses diffusion models to achieve this denoising. It conditionally samples an image generation based on a noisy image as condition, and matches to an intermediate timestep in the diffusion process.



Xiang et al. DDM2: Self-Supervised Diffusion MRI Denoising with Generative Diffusion Models. ICLR 2023.

# AlphaFold 3

- Diffusion-based architecture that predicts the joint 3D structure of complexes including proteins, nucleic acids, small molecules, etc. from sequences



Abramson et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3. Nature 2024.

# Two major classes of vision foundation models

**Representation Learners**

Train on huge amounts of data

Objective: learn powerful generalist feature extractors that can be used for downstream computer vision tasks

- Vision Only (e.g. DINOv2)

- Vision-Language (e.g. CLIP)

**Generative Models**

Train on huge amounts of data

Objective: learn to generate images and/or text

- Text -> Image (e.g. StableDiffusion, DALL-E)

- Text, Image -> Text (e.g. Flamingo, GPT4-V)

# Two major classes of vision foundation models

**Representation Learners**

Train on huge amounts of data

Objective: learn powerful generalist feature extractors that can be used for downstream computer vision tasks

- Vision Only (e.g. DINOv2)

- Vision-Language (e.g. CLIP)

**Generative Models**

Train on huge amounts of data

Objective: learn to generate images and/or text

- Text -> Image (e.g. StableDiffusion, DALL-E)

- Text, Image -> Text (e.g. Flamingo, GPT4-V)

# Autoregressive Language Models



LARGE LANGUAGE MODELS WITH TRANSFORMERS (Daniel Jurafsky & James H. Martin 2024)

# How to Link Vision to LLMs?

Internal Linkage → Vision-Language Models

❖ Architecture
- ○ Integrate Visual Features into Intermediate Layers
- ○ Integrate Visual Features into Input Layer
- ○ Integrate Visual Patches into Input Layer

External Linkage → Vision-Language Agents

# LLaVA: Integrate Visual Features into Input Layer



Language Response $\mathbf{X_a}$

Language Model $f_\phi$    Vicuna

MLP

Projection $\mathbf{W}$

Vision Encoder

$\mathbf{Z_v}$

$\mathbf{H_v}$       $\mathbf{H_q}$

$\mathbf{X_v}$ Image      $\mathbf{X_q}$ Language Instruction

CLIP Vision Encoder
ViT-L/14

Stage 1: Pre-training for Feature Alignment
Stage 2: Fine-tuning End-to-End ($W \; and \; \emptyset$ )

LLaVA (Liu et al. 2023)

# LLaVA Training Recipe, Stage 1: Pre-training

Goal: Align visual features to LLM's word embedding space



Data: Converted from image captioning data
➢  595K image-text pairs filtered from CC3M
➢  convert to instruction-following format

Trainable Parameters
➢  Only $W$

LLaVA (Liu et al. 2023)

# LLaVA Training Recipe, Stage 2: Instruction Tuning (Supervised Finetuning)

Goal: Visual Captioner → Visual Assistant (Follow user instructions)

Language Response $\mathbf{X_a}$

Language Model $f_\phi$

Projection $\mathbf{W}$     $\mathbf{Z_v}$     $\mathbf{H_v}$     $\mathbf{H_q}$

Vision Encoder     $\mathbf{X_v}$ Image     $\mathbf{X_q}$ Language Instruction

Data: leverage ChatGPT/GPT-4 for multimodal instruction-following data collection
➤ 158K language-image instruction-following data
➤ 3 Types: Conversation (Multi-Turn), Detailed description, Complex reasoning

Trainable Parameters
LLaVA (Liu et al. 2023)     ➤     $\{W, \emptyset\}$
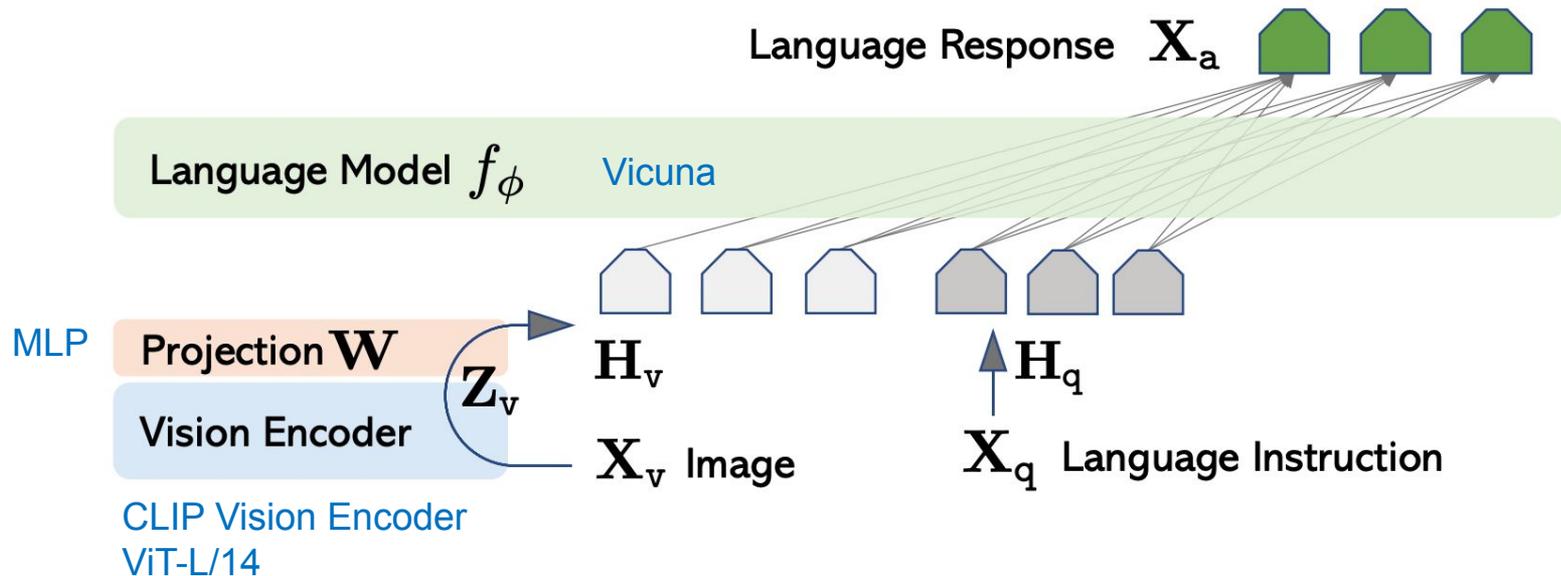
# How to Link Vision to LLMs?

Internal Linkage → Vision-Language Models

❖ Architecture
- ○ Integrate Visual Features into Intermediate Layers
- ○ Integrate Visual Features into Input Layer
- ○ Integrate Visual Patches into Input Layer
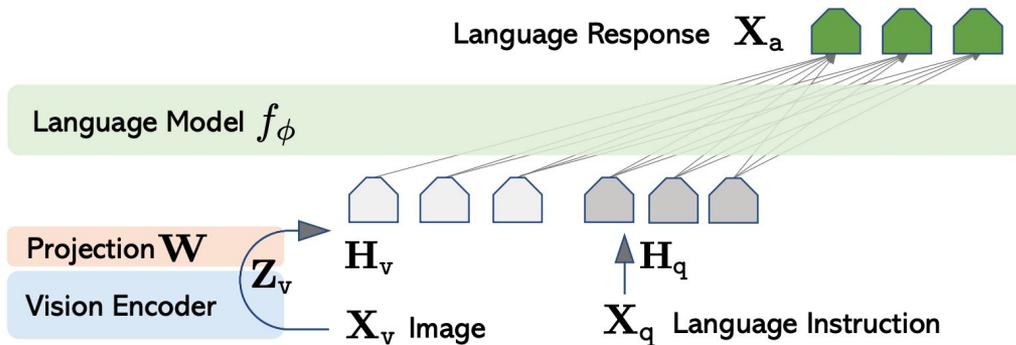
External Linkage → Vision-Language Agents

# Language to Connect Vision Models



HuggingGPT (Shen et al. 2023)

# Code to Connect Vision Models

**Query:** What did the boy do after he dropped the sparkles on the floor?

**Generated code**

```python
def execute_command(video, question, possible_answers):
    video_segment = VideoSegment(video)
    drop_detected = False
    for i, frame in enumerate(video_segment.frame_iterator()):
        if frame.exists("boy") and frame.exists("sparkles") and \
                frame.simple_query("is the boy dropping the sparkles?") == "yes":
            drop_detected = True
            break
    if drop_detected:
        index_frame = i + 1
    else:
        index_frame = video_segment.num_frames // 2
    frame_of_interest = ImagePatch(video_segment, index_frame)
    boys = frame_of_interest.find("boy")
    if len(boys) == 0:
        boys = [frame_of_interest]
    boy = boys[0]
    caption = boy.simple_query("What is this?")
    info = {
        "Caption of frame after dropping the sparkles": caption,
    }
    answer = select_answer(info, question, possible_answers)
    return answer
```
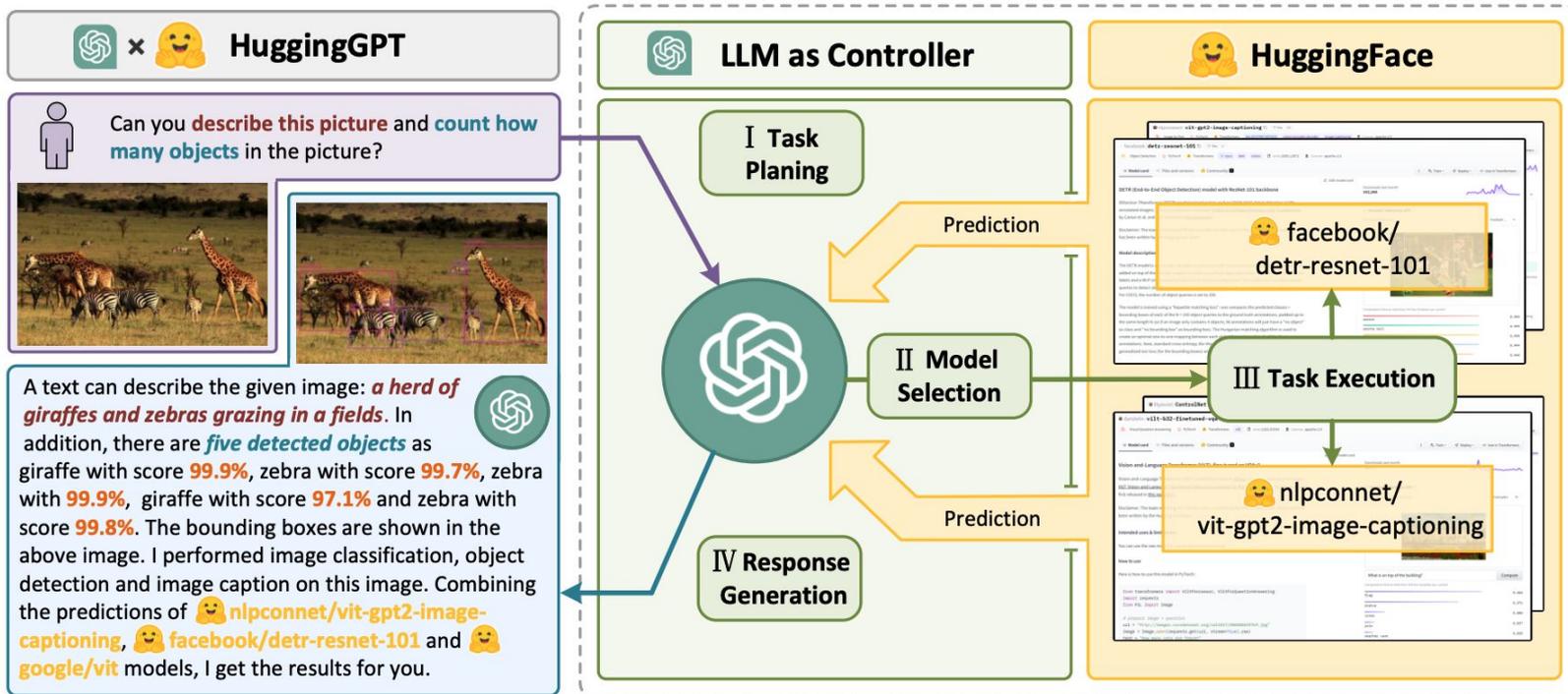
**Execution**

**In:**

```
frame.exists("boy") and \
frame.exists("sparkles") and \
frame.simple_query("is the boy
    dropping the sparkles?") == "yes":
▶ frame = {ImagePatch}
```

▶ i= {int} 25

```
index_frame = i + 1
▶ index_frame = {int} 26
▶ frame_of_interest = {ImagePatch}
```

```
boys = frame_of_interest.find("boy")
▶ boy = {ImagePatch}
```

▶ caption = {str} "a child running with fire in his hands"

▶ answer = {str} "pick it up"

**Result:** *"Pick it up"*

ViperGPT (Surís et al. 2023)

# LLaVA-Med: Training a Large Language-and-Vision Assistant for Biomedicine in One Day

- Extends LLaVA to better answer biomedical questions, using a new biomedical instruction tuning dataset
- Instruction tuning dataset leverages PMC-15M (PubMed Central figure-caption dataset) and covers diverse domains
- Efficiently trained in < 15 hours using eight A100s



(a) Instruction  (b) Responses

Li et al. LLaVA-Med: Training a Large Language-and-Vision Assistant for Biomedicine in One Day. NeurIPS Datasets and Benchmarks 2023.

# Med-Gemini: state-of-the-art generalist biomedical VLM

- Also a generalist (broad domains) biomedical VLM like LLaVA-Med, but extends from the much more powerful Gemini models (and built by internal Google DeepMind team)
- SoTA (state-of-the-art) due to Gemini foundation and additional techniques for medical specialization

Successor to Med-PaLM model based on previous Google PaLM family of models



Saab et al. Capabilities of Gemini Models in Medicine. arXiv 2024.

# LLaVA-Rad: specialized CXR VLM that improves over CheXpert

- Based on LLaVA (similar to LLaVA-Med), but trains from scratch (instead of fine-tuning LLaVA) since it uses a biomedical encoder (Biomed CLIP-CXR) instead of LLaVA's CLIP encoder.

- Trains using 697K image-report pairs from a collection of datasets



Chaves et al. Towards a clinically accessible radiology multimodal model: open-access and lightweight, with automatic evaluation. arXiv 2024.

# PathChat: specialized VLM for pathology

Trained on over 450K instructions from pathology. Will cover further in the discussion presentations!



Lu et al. A multimodal generative AI copilot for human pathology. Nature 2024.

# Two major classes of vision foundation models

## Representation Learners

Train on huge amounts of data

Objective: learn powerful generalist feature extractors that can be used for downstream computer vision tasks

- Vision Only (e.g. DINOv2)

- Vision-Language (e.g. CLIP)

## Generative Models

Train on huge amounts of data

Objective: learn to generate images and/or text

- Text -> Image (e.g. StableDiffusion, DALL-E)

- Text, Image -> Text (e.g. Flamingo, GPT4-V)

# GPU Requirement for Models

VLM  [LLAVA-1.5](#)
- train all models with 8×A100s. 6 hours of pretraining and 20 hours of visual instruction tuning

Vision FM [DINO](#)
- A batch size of 1024, distributed over 16 GPUs. Trained for 3 days.

CLIP [OPEN-CLIP](#)
- ViT-L/14, LAINON-2B data, 384 A100 GPUs for 319 Hours

Vision Generative Model [Stable Diffusion](#)
- using 256 Nvidia A100 GPUs on AWS for a total of 150,000 GPU-hours

# Token Pruning



FastV's Efficiency/Performance Trade-off

FastV could achieve about 45% FLOPs reduction for different LVLMs without sacrificing the performance

An Image is Worth 1/2 Tokens After Layer 2 (Chen et al. 2024)

# Parameter-Efficient Fine-Tuning

☐   Full fine-tuning require more computational cost and becomes infeasible to train on consumer hardware.

☐   Full fine-tuning leads to catastrophic forgetting in the low-data regimes.

☐   Storing and deploying fine-tuned models independently for each downstream task becomes very expensive

   Parameter-Efficient Fine-tuning (PEFT) approaches are meant to address both problems!

❖   Prefix Tuning / Prompt Tuning
❖   Adapter Tuning
❖   LoRA Tuning

# LoRA: Low-Rank Adaptation

**Weight update in regular finetuning**

Outputs

$+$

Pretrained weights $W$

Weight update $\Delta W$

LoRA matrices $A$ and $B$ approximate the weight update matrix $\Delta W$

$d$

Inputs

$d$

**Weight update in LoRA**

Outputs

$+$

Pretrained weights $W$

$B$

$r$

$A$

r can be very small, like 4/8

The inner dimension $r$ is a hyperparameter

Inputs $x$

Random Gaussian initialization for A and zero for B, so $\Delta W$ = BA is zero at the beginning of training

$$h = W_0 x + \Delta W x = W_0 x + B A x$$

LoRA (Hu et al. 2021)

https://magazine.sebastianraschka.com/p/practical-tips-for-finetuning-llms

# The last few class sections: Paper discussions

**Discussion: Recent Advances in Transformers and Vision Representation Learners**

Transformer architecture analysis and improvements
1. [Rotary Position Embedding for Vision Transformer](#)
2. [Vision Transformers Need Registers](#)

Extending SAM vision representation learner to videos
3. [SAM 2: Segment Anything in Images and Videos](#)

Recent example of large-scale vision representation learning for pathology
4. [Virchow2: Scaling Self-Supervised Mixed Magnification Models in Pathology](#)

# The last few class sections: Paper discussions

**Discussion: Recent Advances in Vision-Language Representation Learners**

Techniques for interpretability
1. [Interpreting CLIP's Image Representation via Text-Based Decomposition](#)
2. [Visual Explanations of Image-Text Representations via Multi-Modal Information Bottleneck Attribution](#)

Extending CLIP to multiple modalities
3. [ImageBind: One Embedding Space To Bind Them All](#)

Leveraging CLIP for a unified CT segmentation model across datasets
4. [CLIP-Driven Universal Model for Organ Segmentation and Tumor Detection](#)

# The last few class sections: Paper discussions

**Discussion: Recent Advances Vision Generative Models**

Advancing diffusion model control capabilities
1. [Adding Conditional Control to Text-to-Image Diffusion Models](#)
2. [An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion](#)

Extending to video
3. [Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets](#)

An application connecting RNA sequencing data with image generation
4. [Generation of synthetic whole-slide image tiles of tumours from RNA-sequencing data via cascaded diffusion models](#)

# The last few class sections: Paper discussions

**Discussion: Recent Advances in Vision-Language Generative Models**

Generating both discrete and continuous modalities
1. [Transfusion: Predict the Next Token and Diffuse Images with One Multi-Modal Model](#)

Agentic system for long-form video understanding
2. [VideoAgent: Long-form Video Understanding with Large Language Model as Agent](#)

Instruction tuning on large-scale pathology data
3. [A Multimodal Generative AI Copilot for Human Pathology](#)

# The last few class sections: Paper discussions

**Discussion: Recent Advances in Computing Efficiency**

Token pruning
1. [An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models](#)

IO-aware attention computation
2. [FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness](#)

Parameter-efficient transfer learning from image to video
3. [ST-Adapter: Parameter-Efficient Image-to-Video Transfer Learning](#)

Decomposing medical foundation models into lighter weight expert models
4. [LoRKD: Low-Rank Knowledge Decomposition for Medical Foundation Models](#)

# Where to go from here?

- We hope this class has equipped you with a lay of the land in cutting edge vision and vision-language models, and a foundation for future independent investigation
- Avenues for future exploration may include:
    - **Product-focused projects** for real-world use cases, building more hands-on experience with these models
    - **Research projects** investigating and pushing capabilities of what is possible, building on lines of investigation and insights from recent papers such as those we have covered
    - **Additional courses** on specific model classes (e.g. CS 236 Deep Generative Models), foundation models in healthcare more broadly (less depth on vision, e.g. BIODS 271 Foundation Models for Healthcare), or even broader courses on foundation models (e.g. CS 324 Advances in Foundation Models)

# AMA with the course staff