Topics:

- Generative Adversarial Networks
- LLaVA-OneVision Discussion

CS 4644-DL / 7643-A ZSOLT KIRA

• Assignment 3

- Due July 13th 11:59pm EST
- Projects
 - Project proposal due July 26th

- Next Meta office hours 07/09
 3pm ET on bias/fairness
 - NOT recorded!

W9: July 7

Generative Models (Part I): Generative Adversarial Networks READING: LLaVA-OneVision

W9: July 9 Generative Models (Part II): Diffusion Models PS3/HW3 due July 13th 11:59pm (grace period July 15th)

W10: July 14 Variational Autoencoders (VAEs)

W10: July 16 Open topic! (please suggest). Otherwise default is Object Detection and Segmentation

W11: July 21 Fairness/Bias Discussion, open topics not covered, Wrap-up Final Project Due July 26th 11:59pm (grace period July 28th)

Image to Image CNNs





Classification (Class distribution per image)



Object Detection

(List of bounding boxes with class distribution per box)



Semantic Segmentation (Class distribution per pixel)



Instance Segmentation (Class distribution per pixel with unique ID)





Given an image, output another image

- Each output contains class distribution per pixel
- More generally an image-to-image problem



Semantic Segmentation (Class distribution per pixel)

Instance Segmentation (Class distribution per pixel with unique ID)













Fully connected layers no longer explicitly retain spatial information (though the network can still learn to do so)

Idea: Convert fully connected layer to convolution!

Idea 1: Fully-Convolutional Network





Each kernel has the size of entire input! (output is 1 scalar)

- This is equivalent to Wx+b!
- We have one kernel per output node

Converting FC Layers to Conv Layers



Original:



 $\sum_{k_2=3}^{\infty} k_2 = 3$



Input

Conv Kernel

Output

Larger:







W = 7





Why does this matter?

- We can stride the "fully connected" classifier across larger inputs!
- Convolutions work on arbitrary input sizes (because of striding)

Original sized image



convolutionalization tabby cat heatmap tabby cat heatmap

Larger Image

Larger Output Maps

Long, et al., "Fully Convolutional Networks for Semantic Segmentation", 2015





Convolutional Neural Network (CNN)



Idea 2: "De"Convolution and UnPooling



Example : Max pooling

Stride window across image but perform per-patch max operation

 $X(0:1,0:1) = \begin{bmatrix} 100 & 150 \\ 100 & 200 \end{bmatrix} \longrightarrow max(0:1,0:1) = 200$

Copy value to position chosen as max in encoder, fill reset of this window with zeros





Idea: Remember max elements in encoder! Copy value from equivalent position, rest are zeros



























How can we upsample using convolutions and learnable kernel?

Normal Convolution



Transposed Convolution (also known as "deconvolution", fractionally strided conv) Idea: Take each input pixel, multiply by learnable kernel, "stamp" it on output



"De"Convolution (Transposed Convolution)





Transposed Convolution Example

Georgia Tech



Symmetry in Encoder/Decoder





We can start with a pre-trained trunk/backbone (e.g. network pretrained on ImageNet)!





U-Net

You can have skip connections to bypass bottleneck!



Ronneberger, et al., "U-Net: Convolutional Networks for Biomedical Image Segmentation", 2015



Summary

- Various ways to get image-like outputs, for example to predict segmentations of input images
- Fully convolutional layers essentially apply the striding idea to the output classifiers, supporting arbitrary input sizes
 - (without output size depending on what the input size is)
- We can have various upsampling layers that actually increase the size
- Encoder/decoder architectures are popular ways to leverage these to perform general image-to-image tasks





Generative Models: Introduction











Traditional unsupervised learning methods:



Similar in deep learning, but from neural network/learning perspective





Discriminative vs. Generative Models

- Discriminative models model P(y|x)
 - Example: Model this via neural network, SVM, etc.
- Generative models model P(x)

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





Discriminative vs. Generative Models

- Discriminative models model P(y|x)
 - Example: Model this via neural network, SVM, etc.
- Generative models model P(x)
- We can parameterize our model as $P(x, \theta)$ and use maximum likelihood to optimize the parameters given an unlabeled dataset: $\theta^* = \arg \max \prod_{m=1}^{m} p_{model}(x^{(i)}; \theta)$

$$P^* = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}} \left(\boldsymbol{x}^{(i)}; \boldsymbol{\theta} \right)$$
$$= \arg \max_{\boldsymbol{\theta}} \log \prod_{i=1}^m p_{\text{model}} \left(\boldsymbol{x}^{(i)}; \boldsymbol{\theta} \right)$$
$$= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}} \left(\boldsymbol{x}^{(i)}; \boldsymbol{\theta} \right).$$

- They are called generative because they can otten generate samples
 - Example: Multivariate Gaussian with estimated parameters μ, σ

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks

Generative Models





Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





PixelRNN & PixelCNN





Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





We can use chain rule to decompose the joint distribution

- Factorizes joint distribution into a product of conditional distributions
 - Similar to Bayesian Network (factorizing a joint distribution)
 - Similar to language models!

p

$$(x) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$
Same as language modeling!
$$p(s) = \prod_i p(W_i \mid W_{i-1}, \dots, W_1)$$

$$next$$
word
history
word

- Requires some ordering of variables (edges in a probabilistic graphical model)
- We can estimate this conditional distribution as a neural network

Oord et al., Pixel Recurrent Neural Networks





 Language modeling involves estimating a probability distribution over sequences of words.

$$p(\mathbf{s}) = p(w_1, w_2, \dots, w_n) = \prod_{\substack{i \\ wor}} p(w_i \mid w_{i-1}, \dots, w_1)$$

RNNs are a family of neural architectures for modeling sequences.









$$p(x) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$
$$p(x) = p(x_1) \prod_{i=2}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

1. Choose ordering (upper left, top to bottom, left to right.

Separate out pixel 1

Oord et al., Pixel Recurrent Neural Networks







 $p(x) = p(x_1)p(x_2|x_1)p(x_3|x_1)\prod_{i=1}^{n^2} p(x_i|x_1, \dots, x_{i-1})$

- Model this as RNN with parameters
- Training:
 - We can train similar to language models:
 - Maximum likelihood approach

Downsides?

- Downsides:
 - Slow sequential generation process
 - Only considers few context pixels

Oord et al., Pixel Recurrent Neural Networks



Factorized Models for Images





- Idea: Represent conditional distribution as a convolution layer!
 - Because of spatial locality in images
- Considers larger context (receptive field)
- Practically can be implemented by applying a mask, zeroing out "future" pixels
- Faster training but still slow generation
 - Limited to smaller images

Oord et al., Conditional Image Generation with PixelCNN Decoders





occluded

completions

original



Oord et al., Conditional Image Generation with PixelCNN Decoders

Example Results: Image Completion (PixelRNN)




Geyser



Hartebeest



Grey whale



Tiger

Can we update this to modern times?

Oord et al., Conditional Image Generation with PixelCNN Decoders

Example Images (PixelCNN)





Chameleon: Mixed-Modal Early-Fusion Foundation Models

Multi/Mixed-Modal Large Language Models



Generative Adversarial Networks (GANs)





Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





• *Implicit* generative models do not actually learn an explicit model for p(x)

- Instead, learn to generate samples from p(x)
 - Learn good feature representations
 - Perform data augmentation
 - Learn world models (a simulator!) for reinforcement learning

How?

What architecture lets us generate images? How do we generate a different ^{Generator} image every time?

- Decode architecture
- Learn to sample from a neural network output





• We would like to sample from p(x) using a neural network

Idea:

- Sample from a simple distribution (Gaussian)
- Transform the sample to p(x)



0



Input can be a vector with (independent) Gaussian random numbers

We can use a CNN to generate images!



How do we train this (loss)?

0



Instead, learn to generate samples from p(x)

How?

- Adversarial training that uses one network's predictions to train the other (dynamic loss function!)
- Lots of tricks to make the optimization more stable





- Goal: We would like to generate *realistic* images. How can we drive the network to learn how to do this?
- Idea: Have another network try to distinguish a real image from a generated (fake) image
 - Why? Signal can be used to determine how well it's doing at generation
 - Can be seen as a dynamic (adversarial) loss!













 $N(\mu, \sigma)$

Neural Network

p(x)

Adversarial Networks





Question: What loss functions can we use (for each network)?





Since we have two networks competing, this is a mini-max two player game

- Ties to game theory
- Not clear what (even local) Nash equilibria are for this game

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





Since we have two networks competing, this is a mini-max two player game

- Ties to game theory
- Not clear what (even local) Nash equilibria are for this game
- The full mini-max objective is:

$$\begin{split} \min_{G} \max_{D} V(D,G) &= \mathbb{E}_{\boldsymbol{x} \sim p_{\mathsf{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))] \\ & \text{Sample from real} \end{split}$$

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$$

where D(x) is the discriminator outputs probability ([0,1]) of real image
x is a real image and G(z) is a generated image

The discriminator wants to maximize this:

- D(x) is pushed up (to 1) because x is a real image
- 1 D(G(z)) is also pushed up to 1 (so that D(G(z)) is pushed down to 0)
- In other words, discriminator wants to classify real images as real (1) and fake images as fake (0)

Discriminator Perspective



$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$$

where D(x) is the discriminator outputs probability ([0,1]) of real image
x is a real image and G(z) is a generated image

The generator wants to minimize this:

- First term: G(..) doesn't appear in it!
- 1 D(G(z)) is pushed down to 0 (so that D(G(z)) is pushed up to 1)
- This means that the generator is **fooling** the discriminator, i.e. succeeding at generating images that the discriminator can't discriminate from real

Generator Perspective



Since we have two networks competing, this is a mini-max two player game

- Ties to game theory
- Not clear what (even local) Nash equilibria are for this game
- The full mini-max objective is:

Sample from fake

$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\mathsf{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$$

Generator *minimizes*

How well discriminator does (0 for fake)

• where D(x) is the discriminator outputs probability ([0,1]) of **real** image

• x is a **real image** and G(z) is a **generated** image

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks

Mini-max Two Player Game



Since we have two networks competing, this is a mini-max two player game

Ties to game theory

Not clear what (even local) Nash equilibria are for this game

The full mini-max objective is: Sample from real Sample from fake $\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\mathsf{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$ **Discriminator** *maximizes* How well discriminator How well discriminator does (1 for real) does (0 for fake) where D(x) is the discriminator outputs probability ([0,1]) of **real** image • x is a real image and G(z) is a generated image Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks

Mini-max Two Player Game



Generative Adversarial Networks (GANs)

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k, is a hyperparameter. We used k = 1, the least expensive option, in our experiments.

for number of training iterations do

for k steps do

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D\left(\boldsymbol{x}^{(i)} \right) + \log \left(1 - D\left(G\left(\boldsymbol{z}^{(i)} \right) \right) \right) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Goodfellow, NeurIPS 2016 Generative Adversarial Nets





The generator part of the objective does not have good gradient properties

 $\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\mathsf{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$

- High gradient when D(G(z)) is high (that is, discriminator is wrong)
- We want it to improve when samples are bad (discriminator is right)





Generative Adversarial Networks (GANs)







- Low-resolution images but look decent!
- Last column are nearest neighbor matches in dataset



c)

Early Results



d)



GANs are very difficult to train due to the mini-max objective

Advancements include:

- More stable architectures
- Regularization methods to improve optimization
- Progressive growing/training and scaling

Goodfellow, NeurIPS 2016 Generative Adversarial Nets





Architecture guidelines for stable Deep Convolutional GANs

DCGAN

- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.



Radford et al., Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks



- Training GANs is difficult due to:
 - Minimax objective For example, what if generator learns to memorize training data (no variety) or only generates part of the distribution?
 - Mode collapse Capturing only some modes of distribution
- Several theoretically-motivated regularization methods
 - Simple example: Add noise to real samples!

$$\lambda \cdot \mathbb{E}_{x \sim P_{real}, \delta \sim N_d(0, cI)} \left[\left\| \nabla_{\mathbf{x}} D_{\theta}(x + \delta) \right\| - k \right]^2$$

Kodali et al., On Convergence and Stability of GANs (also known as How to Train your DRAGAN)





Generative Adversarial Nets: Convolutional Architectures

DI9

Samples from the model look much better!

Radford et al, ICLR 2016



Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Generative Adversarial Nets: Convolutional Architectures

Interpolating between random points in latent space



Radford et al, ICLR 2016



Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n



Brock et al., Large Scale GAN Training for High Fidelity Natural Image Synthesis

Example Generated Images - BigGAN





(a) 128×128

(b) 256×256

(c) 512×512

(d)

Figure 4: Samples from our model with truncation threshold 0.5 (a-c) and an example of class leakage in a partially trained model (d).

Brock et al., Large Scale GAN Training for High Fidelity Natural Image Synthesis







https://www.youtube.com/watch?v=PCBTZh41Ris





Generative Adversarial Networks (GANs) can produce amazing images!

Several drawbacks

- High-fidelity generation heavy to train
- Training can be unstable
- No explicit model for distribution
- Larger number of extensions:
 - GANs conditioned on labels or other information
 - Adversarial losses for other applications





Comparison of Methods



Gradually add Gaussian noise and then reverse

LLaVA-OneVision: Easy Visual Task Transfer

Bo Li^{2,♥} Yuanhan Zhang^{2,♥} Dong Guo¹ Renrui Zhang^{3,♥} Feng Li^{4,♥} Hao Zhang^{4,♥} Kaichen Zhang² Peiyuan Zhang² Yanwei Li^{3,♥} Ziwei Liu² Chunyuan Li¹

¹ByteDance ²S-Lab, NTU ³CUHK ⁴HKUST

https://llava-vl.github.io/blog/llava-onevision

Abstract

We present LLaVA-OneVision, a family of open large multimodal models (LMMs) developed by consolidating our insights into data, models, and visual representations in the LLaVA-NeXT blog series. Our experimental results demonstrate that LLaVA-OneVision is the first single model that can simultaneously push the performance boundaries of open LMMs in three important computer vision scenarios: single-image, multi-image, and video scenarios. Importantly, the design of LLaVA-OneVision allows strong transfer learning across different modalities/scenarios, yielding new emerging capabilities. In particular, strong video understanding and cross-scenario capabilities are demonstrated through task transfer from images to videos.

Problem Statement

- What problem does this paper focus on?
 - Is this new or already explored?
 - Is this important?
 - What key applications this is relevant for?
 - What assumptions does this paper make about





Related Work / Situation of Work

• What prior approaches exist to solve this problem?

The SoTA proprietary LMMs, such as GPT-4V [109], GPT-4o [110], Gemini [131] and Claude-3.5 [3], exhibit excellent performance in versertile vision scenarios, including single-image, multi-image and video settings. In the open research community, existing works typically develop models tailored to each individual scenario separately. Specifically, most focus on pushing the performance limits in single-image scenarios [26, 83, 173, 73, 164, 35], only a few recent papers have begun to explore multi-image scenarios [70, 47]. While video LMMs excel in video understanding, they often do so at the expense of image performance [72, 76]. It is rare to have a single open model that reports excellent performance in all three scenarios. LLaVA-OneVision aims to fill this gap by demonstrating state-of-the-art performance across a broad range of tasks, and showcasing interesting emerging capabilities through cross-scenario task transfer and composition.

To the best of our knowledge, LLaVA-NeXT-Interleave [68] is the first attempt to report good performance in all three scenarios, LLaVA-OneVision inherits its training recipe and data for improved





Related Work / Situation of Work

• Built on LLaVA-NeXT

Compared with LLaVA-1.5, LLaVA-NeXT has several improvements:

1. Increasing the input image resolution to 4x more pixels. This allows it to grasp more visual details. It supports three aspect ratios, up to 672x672, 336x1344, 1344x336 resolution.

2. Better visual reasoning and OCR capability with an improved visual instruction tuning data mixture.

3. Better visual conversation for more scenarios, covering different applications. Better world knowledge and logical reasoning.

4 Efficient deployment and inference with SGL and





Related Work / Situation of Work

(1) Dynamic High Resolution

We design our model at high resolution with an aim to **preserve its data efficiency**. When provided with high-resolution images and representations that preserve these details, the model's capacity to perceive intricate details in an image is significantly improved. It reduces the model hallucination that conjectures the imagined visual content when confronted with low-resolution images. Our 'AnyRes' technique is designed to accommodate images of various high resolutions. We employ a grid configuration of $\{2 \times 2, 1 \times \{2, 3, 4\}, \{2, 3, 4\} \times 1\}$, balancing performance efficiency with operational costs. See our updated LLaVA-1.5 technical report for more details.



Illustration of dynamic high resolution scheme: a grid configuration of 2 imes 2



Slides created for CS886 at UWaterloo


Approach and Key Nugget

- What approach does this paper take?
- What is the key "golden nugget" intuition, idea, etc. that leads to approach
 - Lower-level
 - Higher-level?





Contributions

recipe, Please see the detailed development timeline in Section A. In particular, our paper makes the following contributions:

- *Large multimodal models*. We develop LLaVA-OneVision, a family of open large multimodal models (LMMs) that improves the performance boundaries of open LMMs in three important vision settings, including single-image, multi-image, and video scenarios.
- *Emerging Capabilities with Task Transfer*. Our design in modeling and data representations allow task transfer across different scenarios, suggesting a simple approach to yield new emgerging capabilities. In particular, LLaVA-OneVision demonstrate strong video understanding through task transfer from images.
- *Open-source*. To pave the way towards building a general-purpose visual assistant, we release the following assets to the public: the generated multimodal instruction data, the codebase, the model checkpoints, and a visual chat demo.







Architecture



Figure 1: LLaVA-OneVision network architecture. Left: The current model instantiation; Right: the general form of LLaVA architecture in [83], but is extended to support more visual signals.



Slides created for CS886 at UWaterloo



Multi-Resolution

To strike a balance of performance and cost, we observe that the scaling of resolution is more effective than that of token numbers, and recommend an AnyRes strategy with pooling.



(a) Higher AnyRes with Bilinear Interpolation



Slides created for CS886 at UWaterloo



Data

seneral		Doc Chart	
Sing Anguage	tle-imag 3.2M Math/R ^e	se asonine	oca name

General (36.1%)	ALLaVA Inst (70.0K) [16]	AOKVQA (66.2 K)	Cambrian (filtered) (83.1 K)
CLEVR (0.7 K)	COCO Caption (20.0 K)	Hateful Memes (8.5 K)	IconQA (2.5 K)
Image Textualization (99.6 K)	LLaVA-158K (158.0 K)	LLaVA-Witd (train) (54.5 K)	LLaVAR (20.0 K)
OKVQA (9.0 K)	RefCOCO (50.6 K)	ScienceQA (5.0 K)	ShareGPT4o (57.3 K)
Shate GPT4V (91.0 K)	ST-VQA (17.2 K)	TattyQA (9.9 K)	Vision FLAN (186.1 K)
Visual7W (14.4 K)	VisText (10.0 K)	VizWiz (6.6 K)	VQARAD (0.3 K)
VQAr2 (82.8 K)	VSR (2.2 K)	WebSight (10.0 K)	InterGPS (1.3 K)
Doc/Chart/Screen (20.6%)	AI2D (GPT4V) (4.9 K)	AI2D (InternVL) (12.4 K)	Al2D (Original) (3.2 K)
Chart2Text (27.0 K)	ChartQA (18.3 K)	Diagram Image2Text (0.3 K)	Doc-VQA (10.2 K)
DVQA (20.0 K)	FigureQA (1.0 K)	HiTab (2.5 K)	Infographic VQA (4.4 K)
LRV Chart (1.8 K)	RoBUT SQA (8.5 K)	RoBUT WikiSQL (75.0 K)	RoBUT WTQ (38.2 K)
Screen2Words (15.7 K)	TQA (1.4 K)	UReader Caption (91.4 K)	UReader IE (17.3 K)
UReader KG (37.6 K)	UReader QA (252.9 K)	VisualMRC (3.0 K)	
Math/Reasoning (20.1%)	MAVIS MCollect (87.4 K)	MAVIS Data Engine (100.0 K)	Geo170K QA (67.8 K)
Geometry 3K (2.1 K)	GEOS (0.5 K)	Geometry3K (MathV360K) (9.7 K)	GeoMVerse (MathV360K) (9.3 K)
GeoQA+ (MathV360K) (17.2 K)	MapQA (MathV360K) (5.2 K)	CLEVR-Math (5.3 K)	Geo170K Atign (60.3 K)
MathQA (29.8 K)	Super-CLEVR (8.7 K)	TabMWP (45.2 K)	UniGeo (12.0 K)
OQA (72.1 K)	LRV Normal (10.5 K)	RAVEN (2.1 K)	Visual Genome (86.4K)
General OCR (8.9%)	ChromeWriting (8.8 K)	HME100K (74.5 K)	IIIT5K (20 K)
IAM (5.7 K)	K12 Printing (12.8 K)	OCR-VQA (80.0 K)	Rendered Text (10.0 K)
SynthDog-EN (40.1 K)	TextCaps (21.9 K)	TextOCR (25.1 K)	

ge (14.3%) Magpie Pro (I.3 MT) (150.0 K) Magpie Pro (I.3 ST) (150.0 K) Magpie Pro (Qwen2 ST) (150.0 K)

Figure 4: Single-Image 3.2M. A High-Quality Single-Image Dataset Collection. Left: Data Distribution within Each Category. The outer circle shows the distribution of all data categories and the inner circle shows the distribution of data subsets. Right: The detailed quantities of datasets.



Single-Image (31.2%)	Magpie Pro (90.0K)	Vision FLAN (fittered) (55.8K)	Image Textualization (49.8K)
Cauldron (40.2K)	UReader (39.9K)	Share GPT4V (21.0K)	ALLaVA Inst. (21.0K)
Cambrian (fittered GPT4o) (24.9K)	LLAVA-Wild (train) (10.9K)	LAION-GPT4V (&0K)	LLAVA-158K (7.0K)
Ceo170K-QA (6.8K)	Geo170K-Align (6.0K)	Share GPT4o (5.7K)	TabMWP (4.5K)
LLAVAR GPT4 (4.0K)	MapQA (4.3K)	MathQA (3.0K)	TextOCR (GPT4V) (2.5K)
TextCaps (2.2K)	Science QA (1.9K)	FigureQA (1.8K)	GeoQA+ (1.7K)
AI2D (InternVL) (1.2K)	UniGeo (1.2K)	konQA (1.1K)	LRV-Normal (filtered) (1.1K)
TQA (1.0K)	Geometry3K (1.0K)	Super-CLEVR (0.9K)	AI2D (GPT4V) (0.7K)
VizWiz (0.7K)	VQA-AS (0.6K)	CLEVR-Math (0.5K)	PlotQA (0.5K)
GEOS (0.5K)	InfoVQA (0.9K)	PMC-VQA (0.4K)	Geo3K (0.2K)
VQA-RAD (0.2K)	LRV-Chart (0.2K)		
Multi-Image (43.0%)	NLV R (86.4K)	Co-Instruct (50.0K)	ScanNet (49.9K)
RAVEN (35.0K)	IconQA (34.6K)	VIST (26.0K)	ScanQA (25.6K)
ContrastiveCaption (25.2K)	ALFRED (22.6K)	PlintstonesSV (22.3K)	ImageCode (16.6K)
DreamSim (15.9K)	Birds-to-Words (14.3K)	PororoSV (12.3K)	Spot-the-Diff (10.8K)
nuScenes (9.8K)	VISION (9.9K)	WebQA (9.3K)	Recipe QA-VisualCloze (&7K)
RecipeQA-ImageCoherence (8.7K)	TQA (MI) (8.2K)	AESOP (6.9K)	HQ-Edit-Diff (7.0K)
MagicBrush-Diff (6.7 K)	COMICS-Dialogue (5.9K)	MultiV QA (5.0K)	VizWiz (MI) (4.9K)
CLEVR-Change (3.9K)	NextQA (3.9K)	IEdit (3.5K)	Star (3.0K)
DocV QA (MI) (1.9K)	MIT-PropertyCoherence (1.9K)	MIT-StateCoherence (1.9K)	OCR-VQA (MI) (1.9K)
Video (25.9%)	ActivityNet (6.5K)	Charades (23.6K)	Ego4D (0.8K)
No. (04. (0. (10))	CONTRACT OF ON	N	

Figure 5: OneVision 1.6M. A high-quality single-image, multi-image and video dataset collection. Left: Data Distribution within each category. The outer circle shows the distribution of all data categories and the inner circle shows the distribution of data subsets. Right: The detailed quantities of datasets. "MI" means it is the multi-image version dataset proposed by DEMON [69].



Training

- Stage-1: Language-Image Alignment. The goal is to well align the visual features into the word embedding space of LLMs.
- Stage-1.5: High-Quality Knowledge Learning. To strike a balance between compute-efficiency and injecting new knowledge into LMMs, we recommend to consider the high-quality knowledge for LMM learning. The training configuration mirrors the settings used in Stage-2, ensuring consistency and allowing the model to integrate new information seamlessly.
- *Stage-2: Visual Instruction Tuning.* To teach LMM to solve a diverse set of visual task with preferred responces, we organize the instruction data into different groups, described in Section 4.2. The model is scheduled to train on these groups in order.





Training



		Stage-1	Stage-1.5	8	itage-2
				Single-Image	OneVision
Vision	Resolution	384	384×{2×2, 1×{2,3}, {2,3}×1}	384×{{1×1},, {6×6}}	384×{{1×1},, {6×6}}
	#Tokens	729	Max 729×5	Max 729×10	Max 729×10 (See Fig. 3)
Data	Dataset	LCS	Image (Sec. 4.1)	Image (Sec. 4.2)	(Multi)-Image & Video (Sec. 4.2)
	#Samples	558K	4M	3.2M	1.6M
Model	Trainable	Projector	Full Model	Full Model	Full Model
	0.5B LLM	1.8M	0.8B	0.8B	0.8B
	7.6B LLM	20.0M	8.0B	8.0B	8.0B
	72.7B LLM	72.0M	73.2B	73.2B	73.2B
Training	Batch Size LR: ψ_{vision} LR: $\{\theta_{proj}, \phi_{LLM}\}$ Epoch	$512 \\ 1 \times 10^{-3} \\ 1 \times 10^{-3} \\ 1$	$256/5122 \times 10^{-6}1 \times 10^{-5}1$	$256/5122 \times 10^{-6}1 \times 10^{-5}1$	$256/5122 \times 10^{-6}1 \times 10^{-5}1$

Table 1: Detailed configuration for each training stage of the LLaVA-OneVision model. The table outlines the progression of vision parameters, dataset characteristics, model specifications, and training hyperparameters across different stages of the curriculum learning process. We use a global batch size of 512 for the 0.5B model, and 256 for the 7B and 72B models.





Validation

- How do they validate their approach?
 - What data do they use?
 - What baselines do they compare against?





Capability	Benchmark	LLaVA OneVision-0.5B	LLaVA OneVision-7B	LLaVA OneVision-72B	GPT-4V (V-Preview)	GPT-4o
	+AI2D [53]	57.1%	81.4%	85.6%	78.2%	94.2%
Single-Image	ChartQA [101]	61.4%	80.0%	83.7%	78.5%	85.7%
	†DocVQA [103] (test)	70.0%	87.5%	91.3%	88.4%	92.8%
	†InfoVQA [102] (test) Infographic Understanding	41.8%	68.8%	74.9%	-	-
	MathVerse [165] (vision-mini)	17.9%	26.2%	39.1%	32.8%	50.2%
	MathVista [90] (testmini)	34.8%	63.2%	67.5%	49.9%	63.8%
	MMBench [86] (en-dev)	52.1%	80.8%	85.9%	75.0%	-
	Multi-discip MME [28] (cog./perp.)	240/1238	418/1580	579/1682	517/1409	-
	Multi-discip MMStar [19]	37.5%	61.7%	66.1%	57.1%	-
	Multi-discip MMMU [157] (val)	31.4%	48.8%	56.8%	56.8%	69.1%
	College-level Multi-disp MMVet [153]	29.1%	57.5%	63.7%	49.9%	76.2%
	Mutti-discip SeedBench [66] (image)	65.5%	75.4%	78.0%	49.9%	76.2%
	Mutti-discip; Large-scale †ScienceQA [93] High-school Science	67.2%	96.0%	90.3%	75.7%	-
	ImageDC [65]	83.3%	88.2%	91.2%	91.5%	-
	RealworldQA [141]	55.6%	66.3%	71.9%	61.4%	-
	Vibe-Eval [112]	33.8%	51.7%	50.7%	57.9%	63.1%
	MM-Live Bench [161] (2406)	49.9%	77.1%	81.5%		92.4%
	Internet Content Understanding LLaVA-Wilder [65] (small) Reatworld Chat	55.0%	67.8%	72.0%	81.0%	85.9%
	LLaVA-Interleave [68]	33.3%	64.2%	79.9%	60.3%	-
	MuirBench [135]	25.5%	41.8%	54.8%	62.3%	-
fulti-Image	Mantis [47]	39.6%	64.2%	77.6%	62.7%	-
	BLINK [31]	52.1%	48.2%	55.4%	51.1%	-
	Text-rich VQA [84] OCR, Webpage, Ducument	65.0%	80.1%	83.7%	54.5%	•
	ActivityNetQA [155]	50.5%	56.6%	62.3%	57.0%	-
ideo	EgoSchema [98]	26.8%	60.1%	62.0%	-	-
	PerceptionTest [115]	49.2%	57.1%	66.9%	-	-
	SeedBench [66] (video)	44.2%	56.9%	62.1%	60.5%	-
	LongVideoBench [138] (val)	45.8%	56.3%	63.2%	60.7%	66.7%
	Long Video MLV U [170]	50.3%	64.7%	68.0%	49.2%	64.6%
	Long Video Understanding MVBench [71]	45.5%	56.7%	59.4%	43.5%	-
	Multi-discip VideoChatGPT [97]	3.12	3.49	3.62	4.06	-
	Video Conversation VideoMME [29] Multi-discip	44.0%	58.2%	66.2%	59.9%	71.9%

Table 2: Performance comparison to state-of-the-art commercial models with our LLaVA-OneVision models (0.5B to 72B parameters) across diverse evaluation benchmarks spanning multiple modalities. † indicates that the training set has been observed in our data mixture.

Model	AI2D	ChartQA	DocVQA	InfoVQA	MathVerse	MathVista	MMBench	MME	MMMU
	test	test	val/test	val/test	mini-vision	testmini	en-dev	test	val
Qwen-VL-Max [8]	79.3	79.8	-/93.1	-	23.0	51.0	77.6	2281	51.4
Gemini-1.5-Pro [130]	94.4	87.2	-/93.1	-/81.0	-	63.9	-	-	62.2
Claude 3.5 Sonnet [3]	94.7	90.8	-/95.2	49.7	-	67.7	-	-	68.3
GPT-4V [109]	78.2	78.5*	-/88.4	-	32.8	49.9	75.0	517/1409	56.8
GPT-4o [110]	94.2	85.7	-/92.8	-	50.2	63.8	-	-	69.1
Cambrian-34B [133]	79.7	73.8	-/75.5	-	-	53.2	81.4	-	49.7
VILA-34B [77]	-	-	-	-	-	-	82.4	1762	51.9
IXC-2.5-7B [162]	81.5	82.2	-/90.9	-/70.0	20.0	59.6	82.2	2229	42.9
InternVL-2-8B [22]	83.8	83.3	-/91.6	-/74.8	27.5	58.3	81.7	2210	49.3
InternVL-2-26B [22]	84.5	84.9	-/92.9	-/75.9	31.3	59.4	83.4	2260	48.3
LLaVA-OV-0.5B (SI)	54.2	61.0	75.0/71.2	44.8/41.3	17.3	34.6	43.8	272/1217	31.2
LLaVA-OV-0.5B	57.1	61.4	73.7/70.0	46.3/41.8	17.9	34.8	52.1	240/1238	31.4
LLaVA-OV-7B (SI)	81.6	78.8	89.3/86.9	69.9/65.3	26.9	56.1	81.7	483/1626	47.3
LLaVA-OV-7B	81.4	80.0	90.2/87.5	70.7/68.8	26.2	63.2	80.8	418/1580	48.8
LLaVA-OV-72B (SI)	85.1	84.9	93.5/91.8	77.7/74.6	37.7	66.5	86.6	563/1706	57.4
LLaVA-OV-72B	85.6	83.7	93.1/91.3	79.2/74.9	39.1	67.5	85.9	579/1682	56.8

Model	MMVet	MMStar	S-Bench	S-QA	ImageDC	MMLBench	RealWorldQA	Vibe-Eval	LLaVA-W	L-Wilder
	test	test	image	test	test	2024-06	test	test	test	small
Qwen-VL-Max [8]	-	-	-	-	-	-	-	-	-	-
Gemini-1.5-Pro [130]	-	-	-	-	-	85.9	70.4	60.4	-	-
Claude 3.5 Sonnet [3]	75.4	-	-	-	-	92.3	59.9	66.2	102.9	83.1
GPT-4V [109]	49.9	57.1	49.9	75.7	91.5	-	61.4	57.9	98.0	81.0
GPT-40 [110]	76.2	-	76.2	-	92.5	92.4	58.6	63.1	106.1	85.9
Cambrian-34B [133]	-	-	-	85.6	-	-	67.8	-	-	-
VILA-34B [77]	53.0	-	75.8	-	-	-	-	81.3	-	-
IXC-2.5-7B [162]	51.7	59.9	75.4	-	87.5	-	67.8	45.2	78.1	61.4
InternVL-2-8B [22]	60.0	59.4	76.0	97.0	87.1	73.4	64.4	46.7	84.5	62.5
InternVL-2-26B [22]	65.4	60.4	76.8	97.5	91.0	77.2	66.8	51.5	99.6	70.2
LLaVA-OV-0.5B (SI)	26.9	36.3	63.4	67.8	83.0	43.2	53.7	34.9	71.2	51.5
LLaVA-OV-0.5B	29.1	37.5	65.5	67.2	83.3	49.9	55.6	33.8	74.2	55.0
LLaVA-OV-7B (SI)	58.8	60.9	74.8	96.6	85.7	75.8	65.5	47.2	86.9	69.1
LLaVA-OV-7B	57.5	61.7	75.4	96.0	88.9	77.1	66.3	51.7	90.7	67.8
LLaVA-OV-72B (SI)	60.0	65.2	77.6	91.3	91.5	84.4	73.8	46.7	<i>93.</i> 7	72.9
LLaVA-OV-72B	63.7	66.1	78.0	90.3	91.2	81.5	71.9	50.7	93.5	72.0

Table 3: LLaVA-OneVision performance on single-image benchmarks. *GPT-4V reports 4-shot results on ChartQA. All results are reported as 0-shot accuracy.



Strengths / Weaknesses

• Strengths?

• Weakness / Limitations?



Slides created for CS886 at UWaterloo

