

Multimodal LLMs

Yu Meng

University of Virginia

yumeng5@virginia.edu

Nov 11, 2024

Announcement

Join at
slido.com
#2705 251



- Guest lecture (11/08) grades posted; contact Xu (ftp8nr@virginia.edu) if you have questions
- Assignment 5 has been released (deadline: **12/02 11:59pm**)



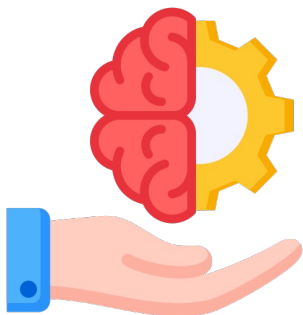
Overview of Course Contents

- Week 1: Logistics & Overview
- Week 2: N-gram Language Models
- Week 3: Word Senses, Semantics & Classic Word Representations
- Week 4: Word Embeddings
- Week 5: Sequence Modeling and Neural Language Models
- Week 6-7: Language Modeling with Transformers (Pretraining + Fine-tuning)
- Week 8: Large Language Models (LLMs) & In-context Learning
- Week 9-10: Reasoning, Knowledge, and Retrieval-Augmented Generation (RAG)
- Week 11: LLM Alignment
- **Week 12: Language Agents**
- Week 13: Recap + Future of NLP
- Week 15 (after Thanksgiving): Project Presentations



(Recap) Overview: Language Model Alignment

- Ensure language models behaviors are aligned with human values and intent
- “HHH” criteria (Askell et al. 2021):
 - **Helpful:** Efficiently perform the task requested by the user
 - **Honest:** Give accurate information & express uncertainty
 - **Harmless:** Avoid offensive/discriminatory/biased outputs



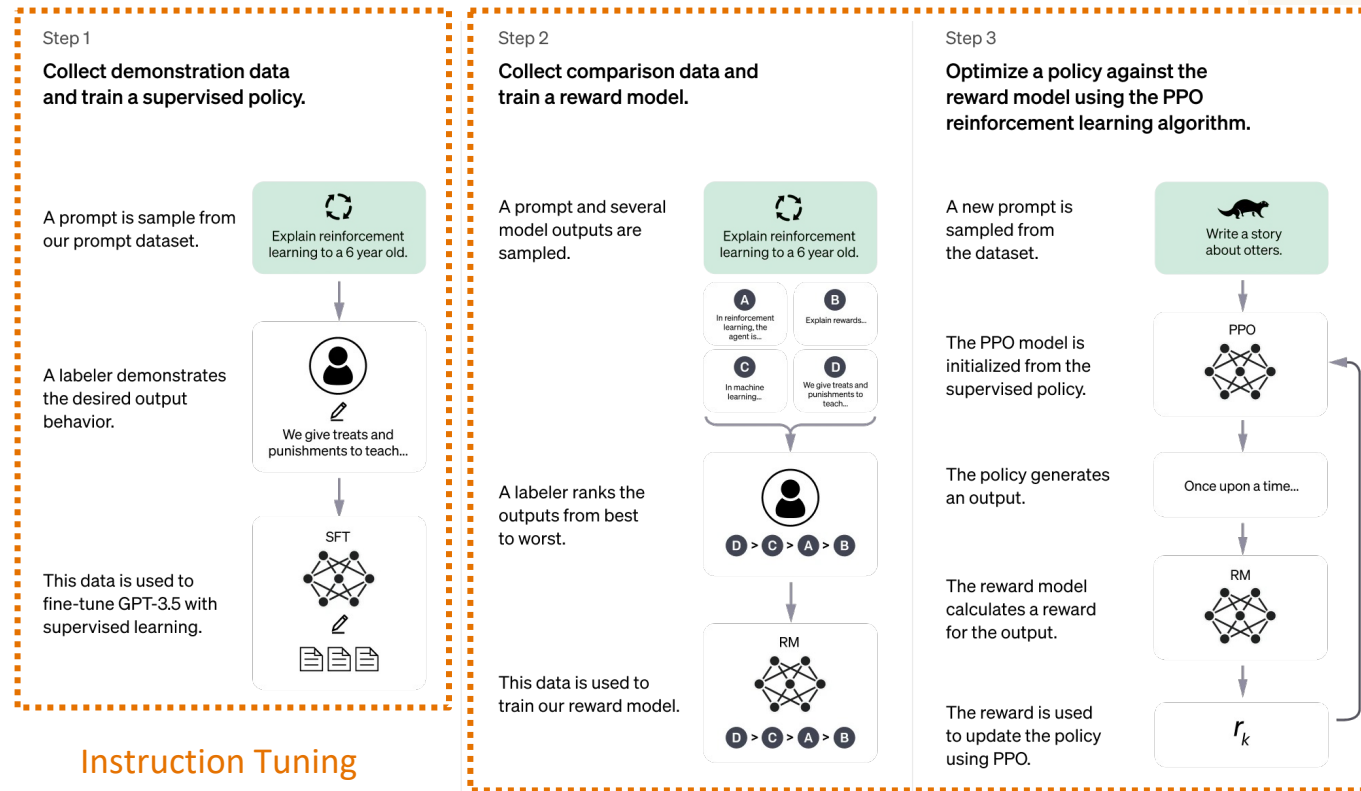


(Recap) Post-training for Alignment

- Pretrained language models are **not** aligned
- Objective mismatch
 - Pretraining is to predict the next word in a sentence
 - Does not involve understanding human intent/values
- Training data bias
 - Text from the internet can contain biased, harmful, or misleading information
 - LMs don't distinguish between good and bad behavior in training data
- (Over-)generalization issues
 - LMs' generalization can lead to outputs that are inappropriate in specific contexts
 - Might not align with intended ethics/honesty standard



(Recap) Language Model Alignment Techniques



Reinforcement Learning from Human Feedback (RLHF)



(Recap) Instruction Tuning: Method

- Input:** task description
- Output:** expected response or solution to the task
- Train LLMs to generate response tokens given prompts

$$\min_{\theta} -\log p_{\theta}(y|x)$$

Response

Prompt

Finetune on many tasks (“instruction-tuning”)

| Input (Commonsense Reasoning) | Input (Translation) |
|---|---|
| Here is a goal: Get a cool sleep on summer days. How would you accomplish this goal? OPTIONS: -Keep stack of pillow cases in fridge. -Keep stack of pillow cases in oven. | Translate this sentence to Spanish: The new office building was built in less than three months. |
| Target keep stack of pillow cases in fridge | Target El nuevo edificio de oficinas se construyó en tres meses. |

Sentiment analysis tasks

Coreference resolution tasks

...

Inference on unseen task type

Input (Natural Language Inference)

Premise: At my age you will probably have learnt one lesson.
Hypothesis: It's not certain how many lessons you'll learn by your thirties.
Does the premise entail the hypothesis?
OPTIONS:
-yes -it is not possible to tell -no

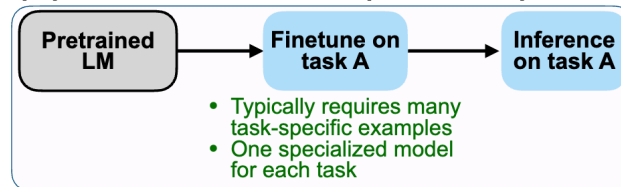
FLAN Response
It is not possible to tell



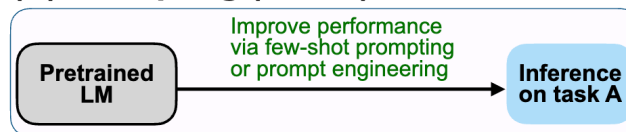
(Recap) Instruction Tuning vs. Other Paradigms

- Task-specific fine-tuning does not enable generalization across multiple tasks
- In-context learning requires few-shot demonstrations
- Instruction tuning enables zero-shot cross task generalization

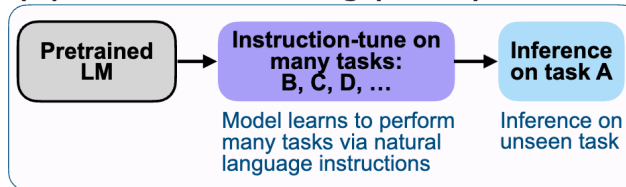
(A) Pretrain–finetune (BERT, T5)



(B) Prompting (GPT-3)



(C) Instruction tuning (FLAN)





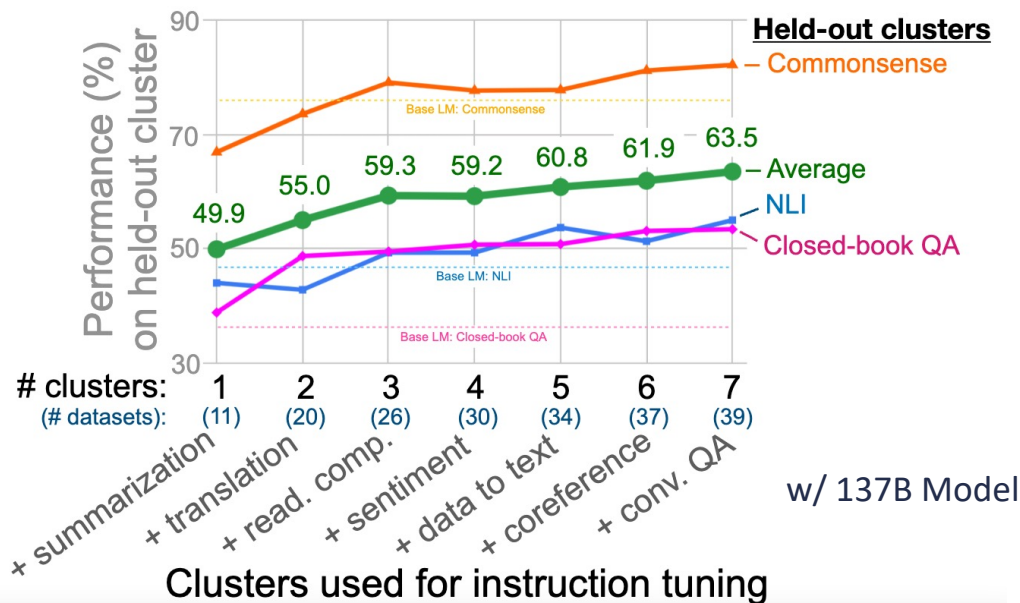
(Recap) Instruction Tuning vs. Pretraining

- Both instruction tuning and pretraining are **multi-task** learning paradigms
- Supervision
 - Pretraining: self-supervised learning (raw data w/o human annotation)
 - Instruction tuning: supervised learning (human annotated responses)
- Task format
 - Pretraining: tasks are implicit (predicting next tokens)
 - Instruction tuning: tasks are explicit (defined using natural language instructions)
- Goal
 - Pretraining: teach LMs a wide range of linguistic patterns & general knowledge
 - Instruction tuning: teach LMs to follow specific instructions and perform a variety of tasks



(Recap) Better Generalization with More Clusters#2705 251

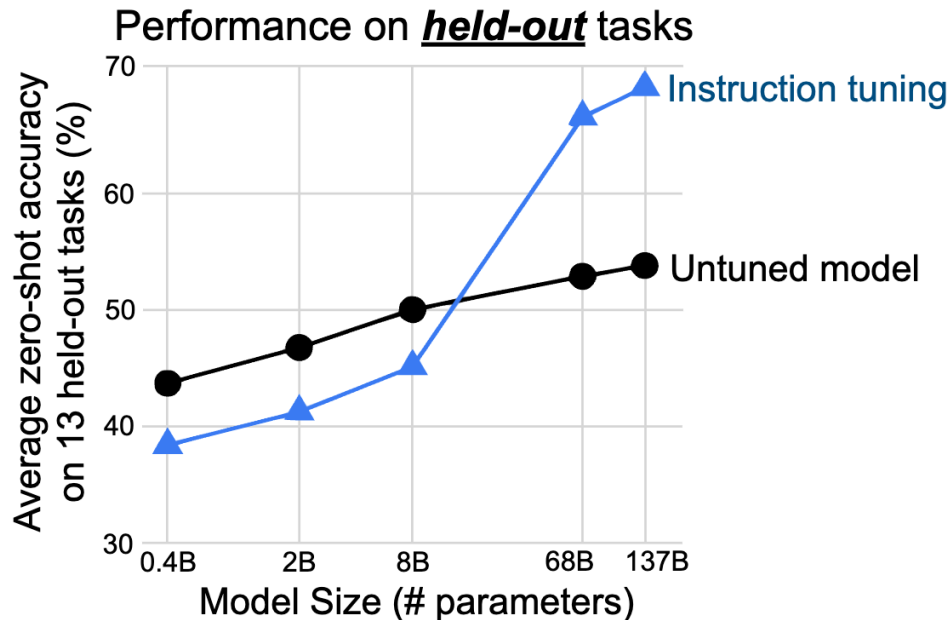
- Held out three clusters from instruction tuning: Commonsense, NLI, Closed-book QA
- More clusters and tasks used in instruction tuning => better generalization to unseen clusters





(Recap) Different Model Sizes

- Instruction tuning can hurt small model ($< 8\text{B}$) generalization
- Instruction tuning substantially improves generalization for large models



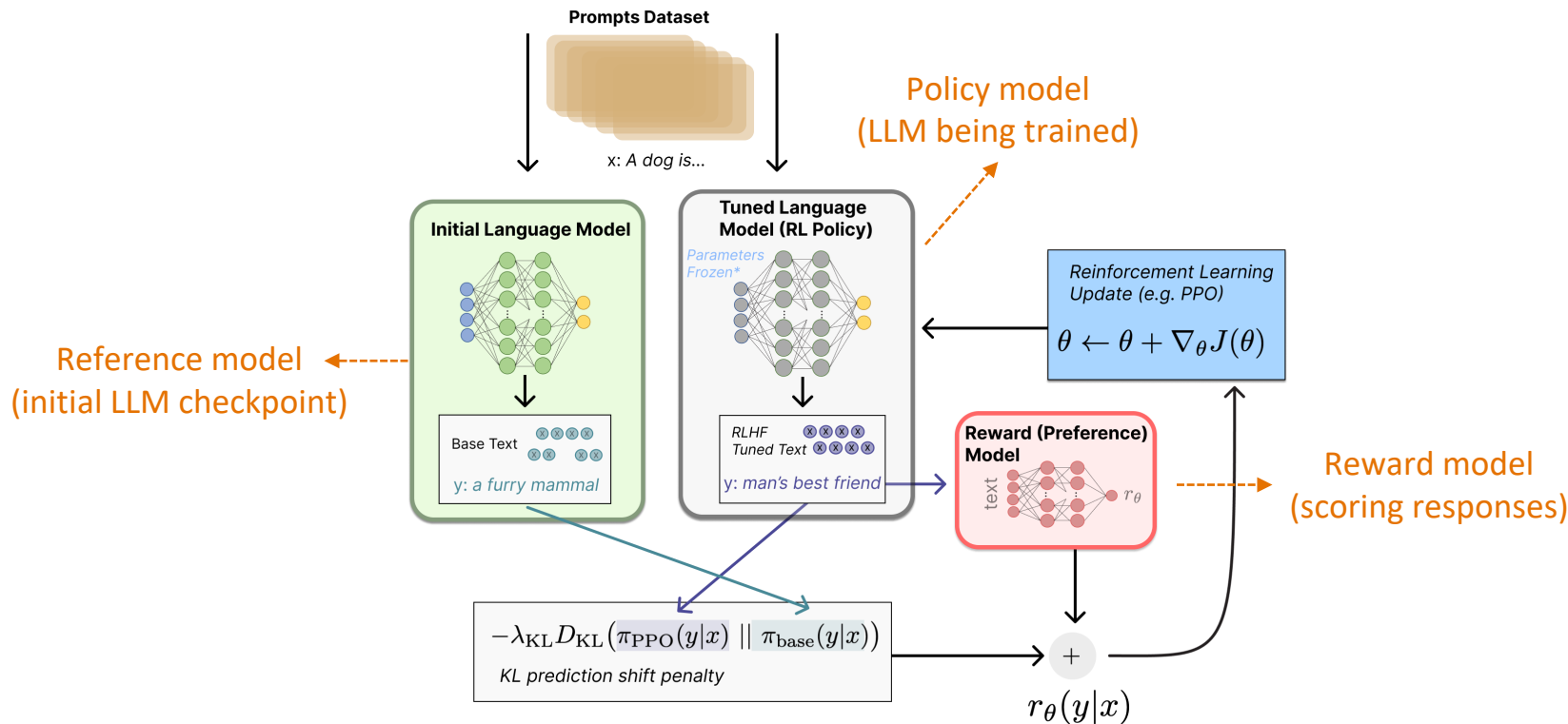


(Recap) Why RLHF

- **Costly human annotations**
 - Instruction tuning requires human annotators to write down the entire expected responses
 - RLHF only relies on preference labels (which response is better?)
- **Open-ended generation**
 - Open-ended creative generation (e.g., story writing) inherently has no single “right” answer
 - RLHF uses human feedback to determine which response is more creative/appealing
- **Token-level learning**
 - Instruction tuning applies the language modeling loss -> penalizes all token mistakes equally regardless of their impact on the overall quality of the output (e.g., a grammatical error might be less critical than a factual inaccuracy)
 - RLHF uses human feedback to prioritize the error types that are more important to correct
- **Suboptimal human answers**
 - Instruction tuning may learn the suboptimal patterns written by humans
 - Identifying a better answer from a few options is usually easier than writing an optimal answer entirely



(Recap) RLHF Illustration





(Recap) Preference Data

Humans are better at relative judgments than absolute ones

What are the steps for making a simple cake?

1. Preheat oven to 350°F (175°C).
2. Grease and flour a cake pan.
3. In a bowl, combine 2 cups flour, 1.5 cups sugar, 3.5 tsp baking powder, and a pinch of salt.
4. Add 1/2 cup butter, 1 cup milk, and 2 tsp vanilla; mix well.
5. Beat in 3 eggs, one at a time.
6. Pour batter into the pan.
7. Bake for 30-35 minutes or until a toothpick comes out clean.
8. Let cool, then frost or serve as desired.

What are the steps for making a simple cake?

1. Warm up the oven.
2. Grease a cake pan.
3. Blend dry ingredients in a bowl.
4. Incorporate butter, milk, and vanilla.
5. Mix in the eggs.
6. Pour into the prepared pan.
7. Bake until golden brown.
8. Add frosting if desired.

Preference data: (x, y_w, y_l)

prompt

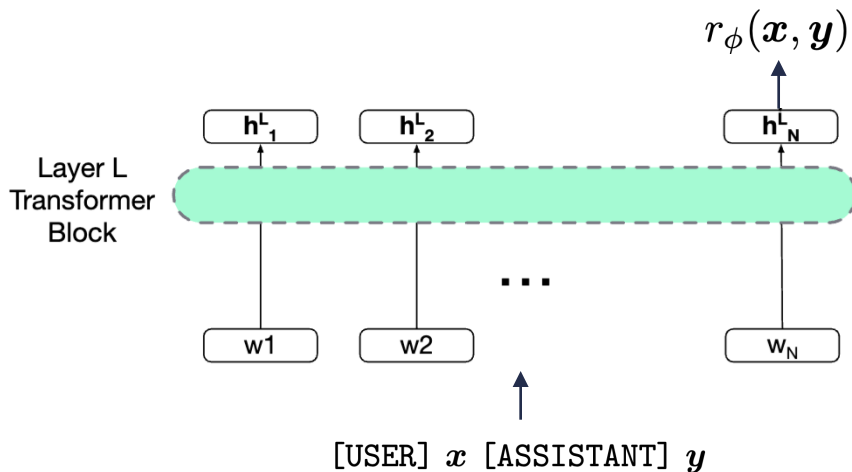
preferred

(winning) response

dispreferred
(losing) response



Goal: train a reward model to assign a higher reward to y_w than y_l



Apply a linear layer at the last token representation to learn a scalar output



(Recap) Reward Model Training

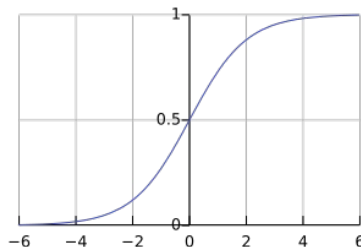
Bradley-Terry pairwise comparison objective

$$\mathcal{L}_{\text{RM}}(r_\phi) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l) \sim \mathcal{D}} [\log \sigma(r_\phi(\mathbf{x}, \mathbf{y}_w) - r_\phi(\mathbf{x}, \mathbf{y}_l))]$$

reward of winning
response

reward of losing
response

$$y = \sigma(x)$$





(Recap) Regularized Reward Optimization

- Add a penalty for drifting too far from the initial SFT checkpoint

$$\max_{\theta} \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\cdot|\mathbf{x})} \left[\underbrace{r_{\phi}(\mathbf{x}, \mathbf{y})}_{\text{Maximize reward}} - \underbrace{\beta \log \left(\frac{p_{\theta}(\mathbf{y}|\mathbf{x})}{p_{\text{SFT}}(\mathbf{y}|\mathbf{x})} \right)}_{\text{Prevent deviation from the initial (SFT) model}} \right]$$

hyperparameter

- Penalize cases where $p_{\theta}(\mathbf{y}|\mathbf{x}) > p_{\text{SFT}}(\mathbf{y}|\mathbf{x})$
- In expectation, it is known as the Kullback-Leibler (KL) divergence $\text{KL}(p_{\theta}(\mathbf{y}|\mathbf{x}) || p_{\text{SFT}}(\mathbf{y}|\mathbf{x}))$



(Recap) Optimization with RL

- Why reinforcement learning:
 - No supervised data available (only a reward model)
 - Encourage the model to explore new possibilities (generations) guided by the reward model
- Optimization: policy gradient methods
 - Optimize the policy (LLM) by adjusting the parameters in the direction that increases expected rewards
- REINFORCE (simplest policy gradient method):

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(a|s) R$$

Diagram illustrating the REINFORCE update rule components:

- α : step size
- ∇_{θ} : policy model (LLM)
- $\log \pi_{\theta}(a|s)$: action (generating the response)
- R : state (user prompt + conversation history)
- R : cumulative reward



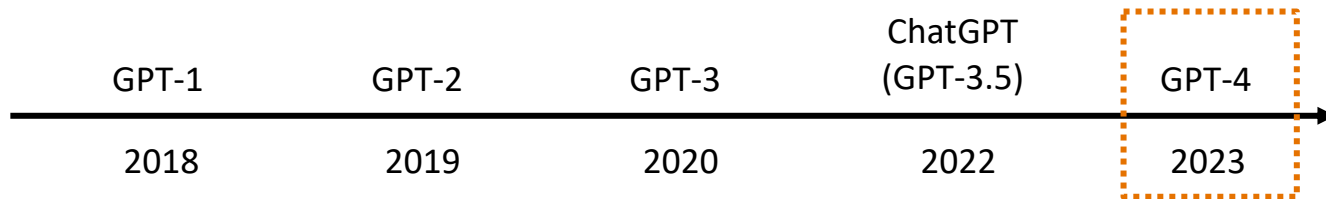
Further Reading on RLHF

- [RAFT: Reward rAnked FineTuning for Generative Foundation Model Alignment](#) [Dong et al., 2023]
- [Iterative Preference Learning from Human Feedback: Bridging Theory and Practice for RLHF under KL-Constraint](#) [Xiong et al., 2023]
- [SLiC-HF: Sequence Likelihood Calibration with Human Feedback](#) [Zhao et al., 2023]
- [SimPO: Simple Preference Optimization with a Reference-Free Reward](#) [Meng et al., 2024]



The Evolution of GPT Models: GPT-4

- GPT-1: decoder-only Transformer pretraining
- GPT-2: language model pretraining is multi-task learning
- GPT-3: scaling up & in-context learning
- ChatGPT: language model alignment
- **GPT-4: multimodality**



GPT-4: Multimodal Input Processing

Join at
slido.com
#2705 251



User What is funny about this image? Describe it panel by panel.



Source: hmmm (Reddit)

accept both images and texts as input

GPT-4 The image shows a package for a "Lightning Cable" adapter with three panels.

Panel 1: A smartphone with a VGA connector (a large, blue, 15-pin connector typically used for computer monitors) plugged into its charging port.

Panel 2: The package for the "Lightning Cable" adapter with a picture of a VGA connector on it.

Panel 3: A close-up of the VGA connector with a small Lightning connector (used for charging iPhones and other Apple devices) at the end.

The humor in this image comes from the absurdity of plugging a large, outdated VGA connector into a small, modern smartphone charging port.

Screenshot source: <https://openai.com/index/gpt-4-research/>



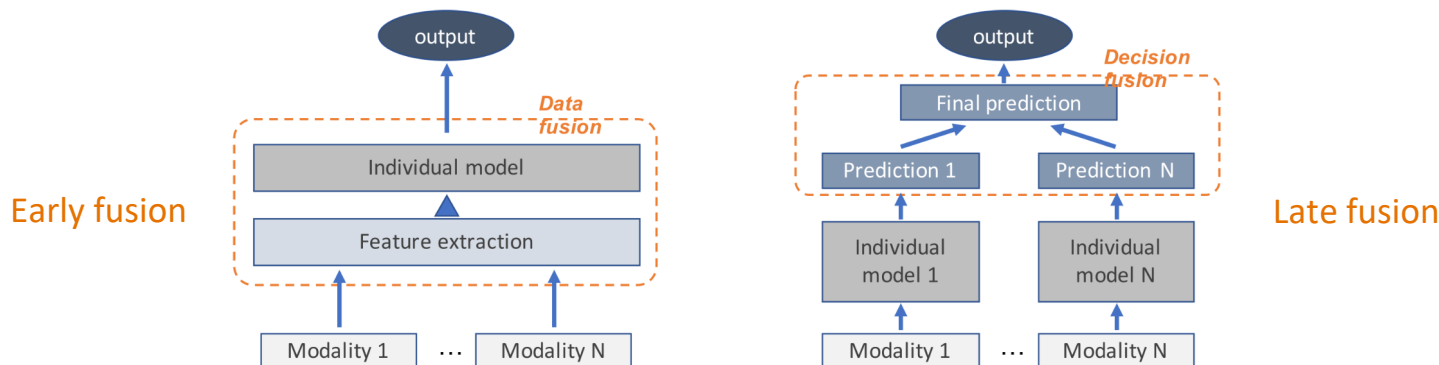
Overview: Multimodal LLMs

- Process and understand multiple types of data (e.g., text, images, audio, and video)
- More comprehensive and contextually rich understanding & generation
- Multimodal input processing (common):
 - Accept and process different types of input data
 - Examples: understanding the content of an image, transcribing and interpreting speech, analyzing video content, or integrating information from sensor data
- Multimodal output generation (less common):
 - Generate output in various modalities
 - Examples: creating realistic images from text descriptions, translating speech to text, or generating music according to user descriptions



Overview: Multimodal Architecture

- Architecture:
 - Require modality-specific architectures (e.g., vision/audio/video encoders)
 - Usually LLMs serve as the strong base
- Multimodal fusion: fuse information from different modalities
 - Early fusion: Combine raw input data from different modalities before processing
 - Late fusion: Process each modality separately and then combine the representations later





Overview: Multimodal Datasets

Training datasets need to contain paired examples of different modalities => teach the model the relationships between different types of data



Source: <https://www.barnorama.com/wp-content/uploads/2016/12/03-Confusing-Pictures.jpg>

What is unusual about this image?

The unusual aspect of this image is a man ironing clothes on the back of a minivan or van. This is not a typical place to perform this activity, as one would usually iron clothes in a more stationary and safe location, such as a home, using a regular ironing board.

Figure source: <https://arxiv.org/pdf/2304.08485>



Learning Aligned Visual Representations

- Goal: learn a joint embedding space where images and their matching text descriptions are close together
- **CLIP** (Contrastive Language-Image Pretraining): predict the correct pairings of a batch of (image, text) training examples

Learning Transferable Visual Models From Natural Language Supervision

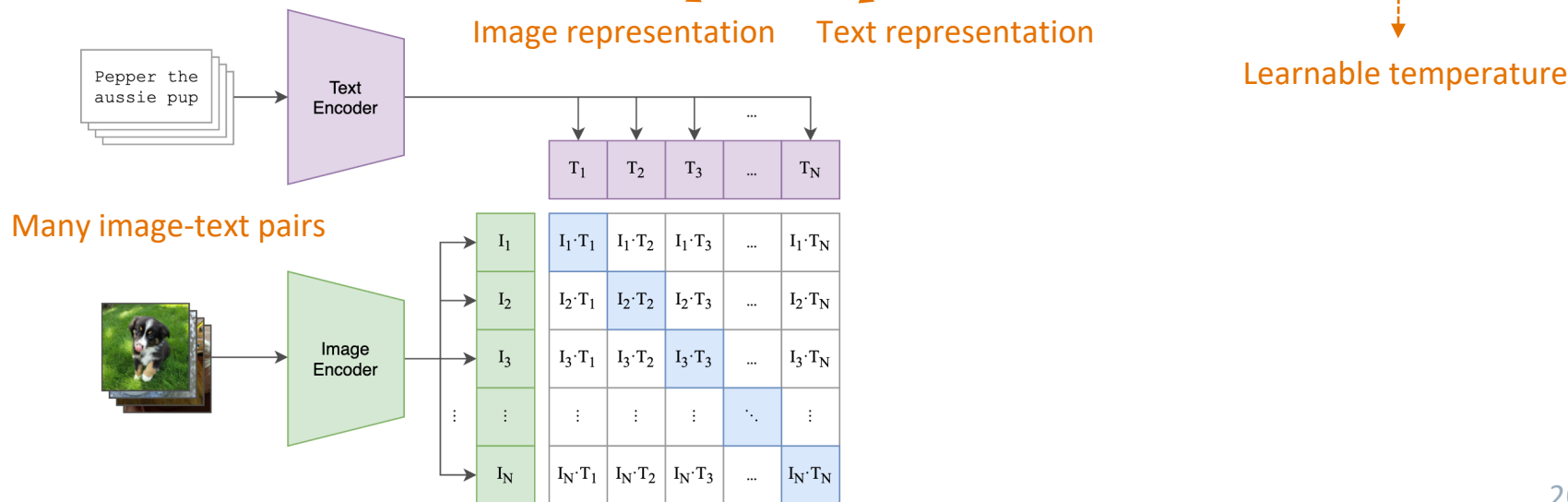
Alec Radford^{*1} Jong Wook Kim^{*1} Chris Hallacy¹ Aditya Ramesh¹ Gabriel Goh¹ Sandhini Agarwal¹
Girish Sastry¹ Amanda Askell¹ Pamela Mishkin¹ Jack Clark¹ Gretchen Krueger¹ Ilya Sutskever¹



CLIP: Contrastive Pretraining

Maximize similarity between correct image-text pairs and minimize for incorrect ones

$$\mathcal{L}_{\text{CLIP}} = -\frac{1}{N} \sum_{i=1}^N \left(\log \frac{\exp(\cos(\mathbf{I}_i, \mathbf{T}_i)/\tau)}{\sum_{j=1}^N \exp(\cos(\mathbf{I}_i, \mathbf{T}_j)/\tau)} + \log \frac{\exp(\cos(\mathbf{T}_i, \mathbf{I}_i)/\tau)}{\sum_{j=1}^N \exp(\cos(\mathbf{T}_i, \mathbf{I}_j)/\tau)} \right)$$

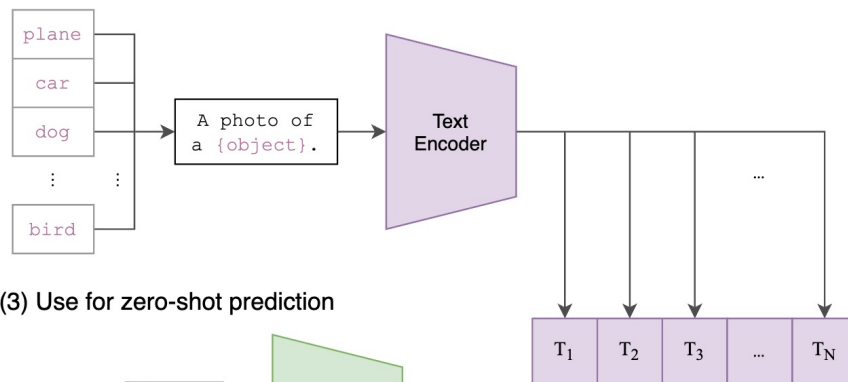




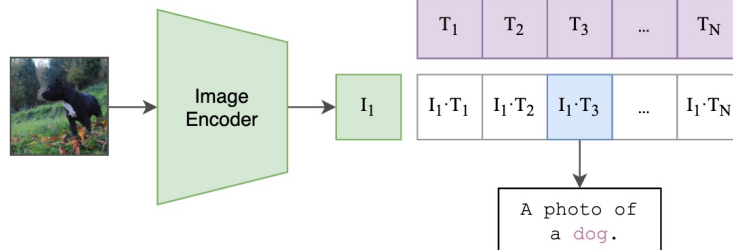
CLIP: Zero-shot Generalization

After training, the text encoder/image encoder can embed the target class names/test images for zero-shot image classification

(2) Create dataset classifier from label text



(3) Use for zero-shot prediction





Visual Instruction Tuning

- Goal: fine-tune a multimodal LLM to learn to follow instructions for tasks that involve both visual and textual information
- **LLaVA** (Large Language and Vision Assistant): combine a pretrained vision encoder (e.g., CLIP) with a large language model (e.g., Llama) for visual instruction tuning

Visual Instruction Tuning

Haotian Liu^{1*}, Chunyuan Li^{2*}, Qingyang Wu³, Yong Jae Lee¹

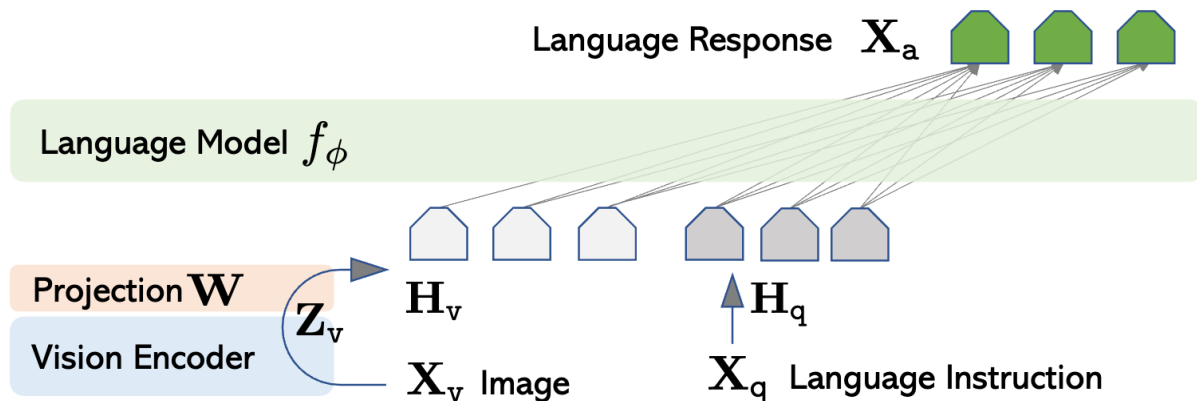
¹University of Wisconsin–Madison ²Microsoft Research ³Columbia University

<https://llava-vl.github.io>



LLaVA: Architecture

- Learn a projection matrix (\mathbf{W}) to convert image representations (\mathbf{Z}_v) to text embeddings (\mathbf{H}_v)
- Concatenate visual tokens (\mathbf{H}_v) with text tokens (\mathbf{H}_q) as input to the model



Adopted in latest
multimodal Llama models

[meta-llama/Llama-3.2-90B-Vision](#)

[meta-llama/Llama-3.2-11B-Vision](#)

LLaVA: Results

Join at

slido.com

#2705 251



Sometimes I just look at pictures of the earth from space and I marvel at how beautiful it all is.



User
LLaVA

Can you explain this meme in detail?

The meme in the image is a creative and humorous take on food, with a focus on chicken nuggets, as the center of the universe. The meme begins with the phrase "Sometimes I just look at pictures of the Earth from space and I marvel at how beautiful it all is..." The main image features a tray of cooked chicken nuggets, with the chicken pieces positioned to resemble the continents and islands on a world map. The meme playfully suggests that the chicken nuggets represent the Earth, and the various locations depicted in the photo are



Further Reading on Multimodal LLMs

- [Zero-Shot Text-to-Image Generation](#) [Ramesh et al., 2021]
- [Flamingo: a Visual Language Model for Few-Shot Learning](#) [Alayrac et al., 2022]
- [AudioLM: a Language Modeling Approach to Audio Generation](#) [Borsos et al., 2022]
- [Movie Gen: A Cast of Media Foundation Models](#) [Polyak et al., 2024]



Thank You!

Yu Meng

University of Virginia

yumeng5@virginia.edu